



(11)

EP 0 724 223 A1

(12)

EUROPEAN PATENT APPLICATION

(43) Date of publication:
31.07.1996 Bulletin 1996/31

(51) Int Cl.⁶: **G06F 17/30**

(21) Application number: 96300437.9

(22) Date of filing: 23.01.1996

(84) Designated Contracting States:
DE FR GB IT SE

(30) Priority: 24.01.1995 US 377152
24.01.1995 US 377881

(71) Applicant: **TANDEM COMPUTERS
INCORPORATED**
Cupertino, CA 95014-3548 (US)

(72) Inventors:
• Carr, Richard W.
Palo Alto, California 94306 (US)

- **Garrard, Brian**
Speldhurst, Kent, TN3 0NS (GB)
- **Mosher, Malcolm, Jr.**
Los Gatos, California 95030 (US)

(74) Representative:
Cross, Rupert Edward Blount et al
BOULT, WADE & TENNANT
27 Furnival Street
London EC4A 1PQ (GB)

(54) **Remote duplicate database facility with database replication support for online line DDL operations**

(57) A local computer system has a local database, application programs that modify the local database, and a transaction manager that stores audit records in a local audit trail reflecting those application program modifications to the local database. A remotely located computer system has a backup database. A remote data duplication facility (RDF) maintains virtual synchronization of the backup database with the local database. The RDF includes an extractor process executed by the local computer system, and a receiver process and a plurality of updater processes executed by the remote computer system. The extractor process extracts audit records from the local audit trail and transmits those records to the receiver process. The receiver distributes the received audit records to a plurality of image trail files in the remote computer system for processing by updater processes, which initiate redo operations of database modifications denoted in at least a subset of the audit records against the backup database. A catalog manager on the local computer system performs online database restructurings while application programs continue to modify the database. The transaction manager stores a Stop Updaters audit record in the local audit trail when each online database restructuring successfully completes. The extractor process transmits the Stop Updaters audit record to the remote computer system and the receiver process moves a copy each received Stop Updaters audit record into all of the image trails. Finally, each updater process stops execution when it reads a Stop Updaters audit record in its as-

signed image trail file.

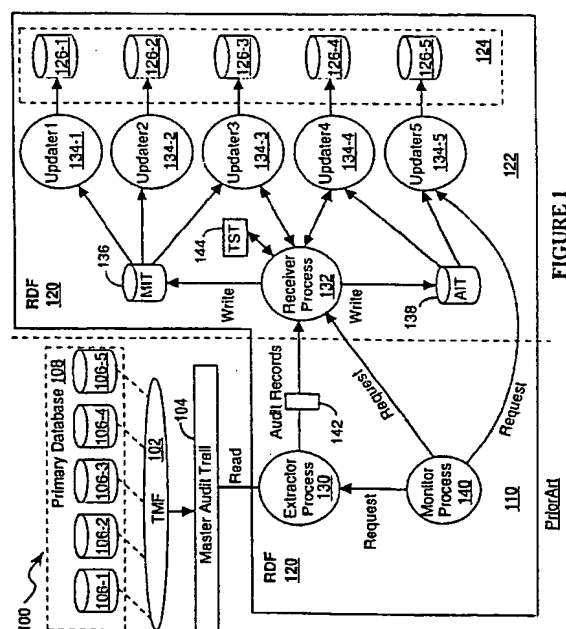


FIGURE 1

Prior Art

Description

The present invention relates generally to database management systems and particularly to a database management system in which a "remote duplicate database facility" monitors changes made to a database on a local system and maintains a copy of that database on a remote system. The present invention also related to a database management system in which database table availability is maintained, with minimal or no user availability outages during table restructuring operations such as splitting a table or index partition, moving an existing table or index partition, creating a new index and moving a table or index partition boundary.

BACKGROUND OF THE INVENTION

The purpose of a "remote duplicate database facility" (hereinafter an "RDF") in a database management system (DBMS) is to provide a higher degree of fault tolerance than can be provided through conventional fault tolerance system architectures involving the use of "shadowed" disk storage (i.e., storing all data to two disk storage systems simultaneously) and parallel computers where one computer takes over if the other fails.

In conventional distributed computer systems designed to perform transaction management, one or more transaction managers (i.e., management processes) generate and store transaction audit entries in an audit trail. Most audit entries denote a database table record event, such as an addition, deletion or modification of a specified database table record in a specified database table. Other audit entries indicate that a specified transaction has committed or aborted. Yet other types of transactions, such as transactions involving restructurings of the database tables (e.g., adding an additional alternate index to a database table or changing the number of partitions for a database table), while still other audit entries may denote other events not relevant to this document.

An RDF system monitors the information added to the audit trail so as to monitor changes made to a database on a primary system (often called the local system), and maintains a backup copy of that database on a remote system by applying the same changes to the database on a remotely located backup system (often called the remote system). In this manner the backup database on the remote system is kept continuously up to date by the RDF with changes made to the local system.

The remote system is preferably located sufficiently far from the primary system that even a widespread disaster, such as loss of power over a large geographic area, that adversely affects the primary system will not affect the remote backup system. The use of an RDF system makes it possible to switch business applications from the primary system to the remote backup system in a short period of time, preferably in a matter of minutes. Additionally, if a planned shutdown of the primary system is necessary, business applications can be stopped and immediately restarted on the backup system to access the replicated database.

An DDL statement is a data definition language statement. DDL statements are used to create and modify database tables. A DDL operation is one that creates or modifies database tables in response to execution of a DDL statement.

Database configuration and reconfiguration operations can have a significant effect on the availability of user applications that need access to databases undergoing structural changes. The Tandem™ NonStop™ SQL/MP relational database management system (DBMS), prior to the present invention, allowed read access, but not write access, to the portions of the database table undergoing the restructuring operation. Furthermore, Tandem's prior art RDF system provided no explicit support for database restructuring operations, other than the fact that any database restructuring operation performed on the primary computer system can also be performed (by operator command) on the remote backup system.

Although most users perform database restructuring operations infrequently, their duration can account for thousands of minutes of application outages per year. A discussion of the cost of application outages appears in the article "An Overview of Non Stop SQL/MP," Ho et al., Tandem Systems Review, July 1994. Moreover, once a database system having an RDF is changed to allow "online" database restructurings, the RDF system must also be changed to avoid desynchronization of the backup database with the primary database.

Prior Art Tandem RDF System

Figures 1 and 2 represent the basic architecture of Tandem Computer's RDF system prior to the present invention. The computer system 100 shown in Figure 1 has a transaction management facility 102 that writes audit entries to a master audit trail (MAT) 104. The audit entries indicate changes made to "audited files" on "RDF protected volumes" 106 of a primary database 108 on a primary system 110. All RDF protected volumes are configured to write all transaction audit records to the MAT 104.

The RDF system 120 includes processes on both the primary (local) computer system 110 and a remote backup computer system 122. The RDF 120 maintains a replicated database 124 by monitoring changes made to "audited files" on "RDF protected volumes" 106 on a primary system and applying those changes to corresponding backup

volumes 126 on the backup computer system 122. An "audited file" (sometimes called an "RDF audited file") is a file for which RDF protection has been enabled, and an "RDF protected volume" is a logical unit of disk storage for which RDF protection has been enabled. Thus, an RDF protected volume may include both audited and unaudited files.

On the primary computer systems 110, an RDF extractor process 130 reads the master audit trail (MAT) 104, which is a log maintained by the transaction management facility (TMF) of all database transactions that affect audited files, and sends any audit records associated with RDF-protected volumes to an RDF receiver process 132 on the backup computer system.

The MAT 104 is stored as a series of files with sequentially numbered file names. The MAT files are all of a fixed size (configurable for each system), such as 64Mbytes. The TMF 102 and Extractor 130 both are programmed to progress automatically (and independently) from one MAT file to the next.

The extractor process 130 appends a timestamp to each audit record that it extracts from the master audit trail 104. The appended timestamp is the timestamp of the last transaction to complete prior to generation of the audit record in the MAT 104. The resulting record is called an audit image record, or image record. The extractor process stores each audit image record in a message buffer 142 having a size of about 28K bytes in the preferred embodiment. The extractor process reads up to 28K bytes of audit records from the MAT 104 at a time, and after each such read operation the resulting message buffer 142 is transmitted to the receiver process 132.

The extractor process 130 waits for an acknowledgment message from the receiver process 132 before continuing its processing of audit records in the MAT 104.

The receiver process 132 writes all audit records received from the extractor to a master Image Trail (MIT) 136 as well as to zero or more auxiliary Image Trails (AITs) 138. The contents of all the image trails 136, 138 are identical. The RDF updater processes 134 on the remote backup system 122 read the audit records from either a master image trail 136 or an auxiliary image trail 138 and apply only audit records associated with committed transactions to the backup database 124. Each RDF-protected volume 106 on the primary computer system 110 has its own updater process 134 on the backup computer system 110 that is responsible for applying audit records to the corresponding backup volume 126 on the backup computer system 110 so as to replicate the audit protected files on that volume. Audit records associated with aborted transactions on the primary system are never applied to the database on the remote backup computer system 122.

The audit image records in each image trail 136, 138 are typically read and processed by two to ten updaters 134. Each updater 134 reads all the audit image records in the corresponding image trail, but utilizes only the audit image records associated with the primary disk volume 106 for which that updater is responsible. Thus, in a system having a large number of RDF protected disk volumes, each updater 134 will utilize only a small fraction of the audit image records in the corresponding image trail 136, 138. For instance, in a system with four image trails (136, 138) and sixteen updaters (four per image trail), only about six percent of the audit records read by each updater (on average) will be relevant to that updater.

The receiver process 132 reads the master image trail 136, processing only the transaction commit/abort records in the audit image trail. From the information in those records it builds a transaction status table (TST) 144 indicating the status of each transaction that has either committed or aborted.

Since the updaters 134 only initiate redo operations on the audit image records for transactions that have committed, the updaters 134 require the transaction status information in the transaction status table 144. To obtain that information, each updater 134 requests transaction status information from the receiver process 132 whenever it reads an audit image record for a database table that the updater is assigned to replicate and for which the transaction status is unknown to the updater.

The receiver process 132, in response to each status request, sends the requesting updater process 134 a message that includes not only the status of the transaction identified in the status request, but also the status of the next hundred or so transactions in the transaction status table that completed after the identified transaction. If the receiver process 132 does not yet know the status of the identified transaction, it does not respond to the status request until it receives a commit/abort record concerning the identified transaction.

When an updater process 134 reaches the end of file of the image trail 136, 138 to which it is assigned, it performs a wait for a preselected amount of time, such as 2 to 10 seconds before attempting to reach more audit image records.

Monitor process 140 accepts user commands for controlling the RDF 120. Monitor process 140 also accepts user requests for status information and requests status information from the various processes in the RDF 120 in order to respond to those user requests.

Referring to Figure 2, the extractor process 130 has a backup extractor process 150 and the receiver process 132 has a backup receiver process 152. The extractor backup process is created by the extractor process 130 and is always resident on a different CPU 160 from the extractor process 130 so as to improve the chances that a hardware failure of the extractor process's CPU will not affect the backup extractor process. Similarly, the receiver backup process 152 is created by the receiver process 132 and is always resident on a different CPU 162 from the receiver process 132.

The extractor backup process 152 is dormant while the primary extractor process 130 remains active. An operating

system procedure 156 periodically checks on the primary extractor process 130 to determine if it is still alive. When the operating system determines that the primary extractor process 130 has failed, the extractor backup process 150 takes over execution of the extractor procedures at a takeover location specified by the last checkpoint performed by the primary extractor process 130 (as will be explained in more detail next) utilizing the last checkpointed information 158 from the primary extractor process to establish the backup extractor's process context.

A "checkpoint" operation is defined in this document to mean the storage of information by one process in a backup process. Thus, checkpoints store information in primary memory, not on durable disk or other secondary memory storage.

The extractor process 130 performs a checkpoint whenever (A) it finishes reading one MAT file and begins reading a next MAT file, (B) whenever the extractor receives acknowledgment from the receiver process 132 that a message buffer has been received, and (C) in a number of special context change circumstances not relevant here. When the extractor process 130 performs a checkpoint, the information 158 transferred by the primary extractor process 130 to its backup process includes:

- a takeover location,
- the extractor's current position in the MAT file; and
- all data structures associated with processing audit records, including partially processed audit records, the last seen commit/abort timestamp, and so on.

The extractor process 130 does not durably store a context record, except as part of an orderly shutdown of the extractor process. This durably stored context record is used at startup to determine where the extractor process 130 should start reading in the MAT 104.

In the prior art Tandem RDF 120, the transferred takeover location can be any point in the extractor process's programming.

In the prior art Tandem RDF 120, the receiver process 132 frequently performs checkpoint operations. More specifically, the receiver process 132 performs a checkpoint (A) upon receipt of a message buffer of audit information from the extractor process (before sending an acknowledgement reply message to the extractor process), (B) upon completion of a no-waited write to the image trail disk files (i.e., upon receipt of notification of successful writes to all the image trails from the associated disk processes), (C) whenever it fills up one set of image files and begins writing to a next set of image files, and (D) in a number of special context change circumstances not relevant here.

Whenever the receiver process 132 performs a checkpoint, the information 164 transferred by the primary receiver process 132 to its backup process 152 includes:

- a takeover location,
- the receiver's current position in the image files (i.e., the current end of file position for each image file, which is also the position for the next write operation to each of the image files); and
- all data structures associated with processing audit records, including the complete message buffer received from the extractor process, partially processed blocks of data not yet written to image disk files, and an updater status table. The updater status table indicates the last reported image trail read position for each updater.

Immediately after the receiver 132 receives a message buffer of audit image records it performs a checkpoint, moving a copy of its context record and all the received data to its backup process. Then it sends an acknowledgment reply message back to the extractor process 130. Next, it writes a copy of the received audit image records into the master image trail 136 and into each auxiliary image trail 138. The master image trail 136 and each auxiliary image trail 138 are a series of disk files having sequentially assigned file names.

The receiver process 132 durably stores its context record every five minutes so as to durably store its current image trail file positions. This durably stored context record is used at startup to determine where the receiver process 132 should start writing data into the image trail files.

The entire context of the backup extractor process 150 is determined by the checkpoint information stored in its address space by the last checkpoint operation performed by the primary extractor process. When the backup extractor process starts up due to failure of the primary extractor process 130, the backup extractor process 150 does not perform any data structure initializations. Rather, it immediately starts execution at the takeover location provided by the last checkpoint and uses the checkpointed data structures as its own data structures.

In a similar manner, the entire context of the receiver process 152 is determined by the checkpoint information stored in its address space by the last checkpoint operation performed by the primary receiver process 132. When the backup receiver process 152 starts up due to failure of the primary receiver process 132, the backup receiver process 152 does not perform any data structure initializations. Rather, it immediately starts execution at the takeover location provided by the last checkpoint and uses the checkpointed data structures as its own data structures.

This method of primary/backup failover used in the prior art Tandem RDF requires (A) that checkpoints transmit all information needed by the corresponding backup process to resume execution immediately, and (B) that checkpoints be performed not only at all significant context changes, but at all points where failure to perform a checkpoint could result in the extractor and receiver becoming desynchronized.

While the RDF system 120 described above has functioned well in commercial use for years, it has been determined, through long experience, that the failover mechanism has a number of shortcomings, and that operating requirements associated with the failover mechanism have greatly limited the throughput (i.e., the number of database updates handled per second) of the RDF system 120. In particular, a primary shortcoming of the prior art failover mechanism is that there remain circumstances, while rare, in which a backup process will not have the information needed, causing the extractor and receiver to become desynchronized.

Two databases, such as the primary database and backup database are said to be "synchronized" if they both have identical contents, in terms of data configured for replication on the backup database. Two databases are said to be "virtually synchronized" if in the course of its normal operation the RDF (i.e., all relevant audit records are being transmitted to the remote backup system) can make up for any difference between the two databases, as in the case of a time lag between data arriving in the MAT file and being transmitted by the RDF to the backup database. When the RDF catches up to the primary system, the two are once again fully synchronized.

In the prior art RDF system 120, there is no coordination between the extractor and receiver processes 130, 132 other than use of message sequence numbers in the message buffers transmitted from the extractor process 130 to the receiver process 132. The prior art RDF system simply assumes that virtual synchronization of the extractor and receiver processes is maintained at all times. Thus, the system does not provide an automatic system and method for the extractor and receiver to resynchronize in the event that desynchronization occurs. For instance, if receiver receives a message buffer with a wrong sequence number, the RDF system 120 would crash.

While such desynchronizations are rare, the mechanisms used by RDF system 120 to avoid desynchronization are complex, use a large percentage of the RDF's entire resources, and are not entirely failure proof.

In the prior art RDF system 120, the maximum number of image trails 136, 138 that the receiver can handle is seven. This limit is based on the receiver's checkpointing requirements, which absorb much of the receiver process's resources, and the I/O overhead associated with writing all audit information to the image trails.

Another shortcoming of the prior art RDF system 120 that limits the RDF system's throughput, is the requirement that the extractor process 130 wait until a message buffer is acknowledged by the receiver process 132 before the extractor process resumes processing more audit records. The RDF system 120 thus uses only one message buffer at a time. This requirement is based on the need to maintain virtual synchronization between the primary and backup systems. However, this "wait until reply" requirement greatly reduces the overall rate at which the extractor process 130 can process audit records.

Yet another shortcoming of the prior art RDF system 120, alluded to above, is that the updaters 134 (actually the disk processes called by the updaters) are inefficiently used in systems with large numbers of updaters (e.g., more than eight updaters) because only a small fraction of the records in the image trail processed by each updater will be relevant to that updater.

It is therefore a primary object of the present invention to provide an improved RDF failover mechanism that imposes significantly less overhead on the RDF system, while providing complete failover protection against all single process failures and all multiple process failures. It is thus an object of the present invention to ensure the integrity of the RDF for any type of failure.

A related object of the present invention is to provide an RDF system with tight synchronization between the extractor and receiver processes and procedures that guarantee, regardless of the cause or type of failure, that the two will resynchronize automatically.

Another related object of the present invention is to provide a primary process to backup process failover mechanism for the extractor and receiver processes that does not rely on frequent checkpointing, thereby avoiding the high overhead associated with such checkpointing and the inherent complexity of such checkpointing.

Another object of the present invention is to provide explicit support in an RDF system for online database restructuring operations performed on the primary system.

SUMMARY OF THE INVENTION

In summary, the present invention is a distributed computer database system having a local computer system and a remote computer system. The local computer system has a local database stored on local memory media, application programs that modify the local database, and a transaction manager that stores audit records in a local audit trail reflecting those application program modifications to the local database as well as commit/abort records indicating which of the transactions making those database modifications committed and which aborted. Each audit record has an associated audit trail position in the local audit trail.

The remote computer system, remotely located from the local computer system, has a backup database stored on remote memory media associated with the remote computer system.

A remote data duplication facility (RDF) is partially located in the local computer system and partially in the remote computer for maintaining virtual synchronization of the backup database with the local database. The RDF includes an extractor process executed by said local computer system, and a receiver process and a plurality of updater processes executed by the remote computer system.

The extractor process extracts audit records from the local audit trail. It has a plurality of message buffers for buffering groups of the extracted audit records together and transmits each message buffer to said remote computer system when the buffer is full of a timeout occurs. Each transmitted message buffer has an associated sequence number. The sequence numbers for sequentially transmitted message buffers follow a predefined sequence. The extractor process continues, after transmitting a any message buffer to the remote computer system, to buffer groups of extracted audit records in other ones of the message buffers and to transmit those message buffers to the remote computer system.

The receiver process stores an expected next message sequence number and a context record denoting a restart audit trail position value. The receiver process receives message buffers transmitted by the extractor process. It compares the message sequence number associated with each received message buffer with its locally stored expected next sequence number. If the two do not match, the receiver process transmits an error message to the extractor process. If the two do match, the receiver process responds by sending a reply message to the extractor process acknowledging receipt of the message buffer, updating the expected next message sequence number in accordance with the predefined sequence. It then distributes the audit records in the received message buffer to a plurality of image trail files in the remote computer system and updates the restart audit trail position value based on the audit trail positions associated with said audit records in said received message buffer.

Each updater process, executed by the remote computer system, reads the audit records in an assigned one of said image trail files and initiates redo operations of database modifications denoted in at least a subset of the read audit records against the backup database.

The extractor process responds to each reply message acknowledging receipt of a message buffer by enabling reuse of that message buffer, responds to each error message by determining the restart audit trail position value stored by the receiver process and then extracting audit records from the local audit trail starting at that restart audit trail position value.

The local computer system includes a catalog manager for performing online database restructurings while application programs continue to modify the database. The transaction manager stores a Stop Updaters audit record in the local audit trail when each online database restructuring successfully completes.

The extractor process transmits the Stop Updaters audit record to the remote computer system and the receiver process moves a copy each received Stop Updaters audit record into all of the image trails. Finally, each updater process stops execution when it reads a Stop Updaters audit record in its assigned image trail file.

BRIEF DESCRIPTION OF THE DRAWINGS

Additional objects and features of the invention will be more readily apparent from the following detailed description and appended claims when taken in conjunction with the drawings, in which:

Figure 1 is a block diagram of a prior art database management system with a remote duplicate database facility. Figure 2 is a conceptual representation of the checkpoint and failover procedures used by the system shown in Figure 1.

Figure 3 is a block diagram of a database management system with a remote duplicate database facility in accordance with the embodiment of the present invention.

Figure 4 is a conceptual representation of the checkpoint, context save, and failover procedures used by the system shown in Figure 3.

Figures 5A and 5B depict data structures used by the extractor process in a preferred embodiment of the present invention.

Figures 6A-6E are flowcharts of procedures executed by the extractor process in a preferred embodiment of the present invention.

Figure 7A is a block diagram of a context record for the receiver process in a preferred embodiment of the present invention. Figures 7B-7D are block diagrams of data structures used by the receiver process in a preferred embodiment of the present invention.

Figures 8A-8G flowcharts of procedures executed by the receiver process in a preferred embodiment of the present invention.

Figure 9 is a block diagram of data structures, stored in primary memory, used by each updater process in a preferred embodiment of the present invention.

Figures 10A-10E are flowcharts of procedures executed by the updater processes in a preferred embodiment of the present invention.

Figure 11A is a block diagram of a database table. Figure 11B is a block diagram of an alternate index. Figure 11C is a block diagram of the data structure of an Audit Trail.

Figure 12 is a conceptual diagram of the three primary phases of the database restructuring procedures in embodiments of the present invention.

Figure 13 is a block diagram of a context record for the monitor process in a preferred embodiment of the present invention.

Figure 14 is a flowchart of the monitor procedure for responding to a Stop Updaters message from the extractor process in a preferred embodiment of the present invention.

Figure 15 is a flowchart of the monitor procedure for responding to a Start Update command in a preferred embodiment of the present invention.

Figure 16A is a flowchart of the receiver procedure for responding to a Expect Stop Updaters message from the monitor process in a preferred embodiment of the present invention.

Figure 16B is a flowchart of the receiver procedure for responding to an Updater Close message from an updater process in a preferred embodiment of the present invention.

DESCRIPTION OF THE PREFERRED EMBODIMENTS

Figures 3 and 4 represent the basic architecture of a computer system 200 utilizing the remote duplicate database facility (RDF) 220 of the present invention. Operation of the present invention will be explained by first giving an overview of how the present invention differs from the prior art Tandem RDF system discussed in the background section of this document. The overview is then followed by a detailed explanation of the extractor, receiver and updater processes provided by embodiments of the present invention, and their failover procedures.

Overview of Improved RDF System

The extractor-receiver-updater architecture of the RDF system 220 (shown in Figure 3) has many similarities to the RDF system 120 of Figures 1 and 2. Furthermore, the operation of the transaction management / massively parallel system (TM/MP) 202 is similar to that of the TMF 102, although most of the changes are not directly relevant to the present invention, with the exception of a new type of audit record stored in the MAT 204 when an "online DDL" operation is performed.

The present extractor process 230 differs from the prior art version as follows. First, the present process 230 does not perform frequent checkpoint operations. In fact, the extractor process performs only a single checkpoint operation during startup of the extractor process, and that checkpoint 258 only sends a takeover location to the backup extractor process 250. (See Figure 4.) After that, the extractor process 230 performs no further checkpoints. It also does not durably store a context record. Rather, the extractor process 230 has been revised so that the extractor relies on information received from the receiver process 232 when recovering from a failover, as will be explained in more detail below, as well as during an RDF startup.

The second most significant change to the extractor process is that it now uses a plurality of message buffers 242. The extractor process 230 now uses two to eight message buffers 242, with four message buffers being a typical configuration. After filling and transmitting a message buffer 242 to the receiver process via a communication channel 244, the extractor process 230 does not wait for an acknowledgement reply message from the receiver process 232. Rather, as long another message buffer is available, it continues processing audit records in the MAT 204, storing audit image records in the next available message buffer 242. Each message buffer 242 is made unavailable after it is transmitted to the receiver process 232 until a corresponding acknowledgement reply message is received from the receiver process 232, at which point the message buffer 242 becomes available for use by the extractor process 230.

These two changes (i.e., almost complete elimination of checkpointing and use of multiple message buffers) to the extractor process 230 greatly increase its efficiency, the first by almost eliminating resource usage for failover protection, and the second by allowing the extractor process 230 to continue processing audit records in the MAT 204 while waiting for the receiver process to acknowledge receipt of previously sent audit records.

The present receiver process 232 differs from the prior art version as follows. First, the present receiver process 232 does not perform frequent checkpoint operations. In fact, the receiver process performs only a single checkpoint operation during startup of the receiver process, and that checkpoint 264 only sends a takeover location to the backup receiver process 252. (See Figure 4.) After that, the receiver process 232 performs no further checkpoints. However, it does periodically (e.g., once every 5 to 25 seconds) durably store a context record 270 on a nonvolatile (disk) storage device 272. The context record 270 stored by the receiver process 232 is quite small, consisting primarily of two location values per auxiliary image trail 238 and three location values for the master image trail 236.

Compared with the checkpointing performed by the above described prior art receiver process 232, which was typically performed multiple times per second during normal usage and involved the storage of large quantities of information (up to 30Kbytes of data), the periodic context record save operations by the present receiver process 232 use significantly less resources.

The second significant change to the receiver process is that the receiver process immediately acknowledges each received message buffer. No processing of the message buffer is performed before the acknowledgement is sent. The elimination of message buffer checkpointing by the receiver, and the use of multiple message buffers by the extractor, makes both the receiver and extractor more efficient because idle time in both processes is greatly reduced. Throughput of records from extractor to receiver has improved by approximately a factor of ten (from approximately 110K bytes/sec to approximately 1300K bytes/sec).

A third significant change to the receiver process is that it now sorts received audit records such that (A) commit/abort records are stored only in the master audit trail 236, and (B) each database update audit record is moved into only the one audit trail 236, 238 corresponding to the only updater process 234 which will potentially use that audit record to update data stored on a backup volume 126. Furthermore, there is no limit on the number of image trails used, while there was a limit of seven image trails in the prior art RDF system due to checkpointing limitations. These changes significantly reduce the I/O burden on the receiver process 232. In a system with N audit trails 236, 238, the number of audit records written to disk files is reduced by a factor of N in comparison with the prior art receiver process. This change also makes the updater processes 234 much more efficient than the prior art updater processes 134 because the number of audit records read (by their corresponding disk processes) is also reduced on average by a factor of N.

A fourth significant change to the receiver process is that the receiver process now includes a separate buffer 274 for each image trail. In the preferred embodiment, each image trail buffer is 32K bytes long. Each image trail buffer 274 is structured as seven blocks of length 4K bytes, plus an overflow area and each image trail buffer is considered to be full when the first 28K bytes of the buffer have been filled with audit records.

A fifth change to the receiver process is that whenever it receives a special "Stop Updaters" audit record, it copies that record into all the audit trails. The Stop Updaters audit record, produced on the primary system 210 by special "online DDL" procedures, causes all the Updaters 234 to stop and prompts the operator of the RDF to (A) perform the same DDL procedure on the remote backup system as was performed by the online DDL procedure and then (B) to re-start the updaters. This last change to the receiver process affects neither the failover effectiveness of the RDF system nor the computational efficiency of the receiver process 232. Rather, it is a procedure used to ensure continued virtual synchronization of the local and remote database when "online DDL" procedures are used to restructure database objects with minimal interruption of user access to the database objects being restructured.

Unlike the prior art RDF system, which had virtually no synchronization of extractor and receiver processes, the present system provides tight synchronization of the extractor and receiver processes and provides for automatic resynchronization whenever either process is started or has a failover, and whenever the receiver process receives audit records out of order from the extractor process.

A change to the updater processes is that they now read a master image trail position value embedded periodically in "header" records in the their respective image trail files. Whenever the updater process sends a request to the receiver process, it reports to the receiver process the saved MIT position and the saved image trail file position in the updater's last durably stored context record.

The first reported value is used by the receiver process 232 to determine where in the master image trail 236 to start when reconstructing the transaction status table 144, and also to determine which portions of the transaction status table 144 are no longer needed and thus can be discarded. The reported MIT position is also used to determine which MIT files can be discarded. The second reported value is used by the receiver process 232 to determine which image trail files have been processed by all the updaters assigned thereto and thus can be deleted (as was also the case in the prior art RDF system 120).

Detailed Explanation of Extractor Process

Referring to Figures 5A and 5B, the primary data structures used by the extractor process 230 are as follows. As stated earlier, the extractor process 230 utilizes two or more message buffers 242. A portion of each message buffer 242 is used to store a "header" 280, which includes (A) a message sequence number and (B) a timestamp. The body 282 of the message buffer 242 is used to store audit image records 284. Each image record 284 includes an audit information portion 286, a MAT position value 288 and a timestamp value 290. The audit information portion 286 and the field occupied by the MAT value 288 are copied from an audit record in the MAT 204, while the timestamp field 290 is appended by the extractor process to create an "audit image record" 284.

The audit information portion 286 consists of the standard information found in audit records in the MAT 204, such as before and after field values for a modified row in a database table, or a commit/abort indication for a completed

transaction. In embodiments of the present invention, a new audit image record is defined in which the audit information portion 286 contains a "Stop Updaters" value, indicating that each updater process 234 should be stopped when it reads that audit record.

The extractor process 230 also maintains a message buffer status table 294, which indicates for each message buffer whether that buffer is available for use, not available for use, or is currently in use by the extractor. In addition, the extractor process 230 maintains a message sequence number in register 295, a MAT file pointer in register 296, a local timestamp value in register 297, and a scratch pad 298 in which it stores audit image records that it is currently processing.

Finally, the extractor process 230 includes a data structure 299 for storing reply messages received from the receiver process 232. This data structure includes a first field indicating the type of message received, which is equal to either "message buffer acknowledgement" or "resynch reply", a message buffer identifier, and a "message value" field. The message value field is equal to a MAT position value when the message type is "resynch reply," and is equal to either an "OK" or "Error" condition code when the message type is "message buffer acknowledgement."

Appendix 1 lists a pseudocode representation of the procedures executed by the extractor process 230. Appendix 2 lists a pseudocode representation of the procedures executed by the receiver process 232. Appendix 3 lists a pseudocode representation of the procedures executed by the updater processes 234. The pseudocode used in Appendices 1-3 is, essentially, a computer language using universal computer language conventions. While the pseudocode employed here has been invented solely for the purposes of this description, it is designed to be easily understandable by any computer programmer skilled in the art.

Referring to Figures 6A-6E, and the pseudocode in Appendix 1, the extractor process 230 works as follows.

The Extractor Startup Procedure 300 is called whenever the extractor process 230 or its backup starts up, as in the case of a failover or a transfer of control back to the primary extractor process 230 from the backup extractor process. The Startup procedure begins by creating a backup process (302). The startup procedure then performs a "static initialization" of the extractor process (304), which means that all static data structures used by the extractor process are allocated and initialized. While initializing static data structures, the extractor process reads information denoting the set of RDF protected objects from a disk file having a predefined filename and disk location, and builds an internal table of RDF protected disk volumes. This table is used later as a audit record filter, such that object update audit records for non-RDF protected volumes are ignored by the extractor process. Then a checkpoint operation is performed in which a takeover location is transmitted to the backup extractor process (306). The takeover location is, in essence a program address, and in the preferred embodiment the takeover location is the program location at which execution of the volatile initialization procedure 310 begins. Finally, the Extractor Startup procedure calls (308) the Extractor Volatile Initialization procedure 310.

The Extractor Volatile Initialization procedure 310 is called during startup by the Extractor Startup procedure 300 and when the extractor receives an Error reply message in response to a message buffer. The Extractor Volatile Initialization procedure begins by allocating and initializing all volatile data structures used by the Extractor process, including message buffers 242, the message buffer status array 295 (312), and the message sequence number (which gets initialized to an initial value such as 1. Then the Extractor Volatile Initialization procedure transmits a Resynchronization Request message to the receiver process (314) and waits for a Resynch Reply message (316). The Resynch Reply message will contain a MAT position value, which the Extractor Volatile Initialization procedure moves (318) into the MAT position pointer MATptr 296. Finally, the Extractor Volatile Initialization procedure calls (320) the main Extractor procedure 330.

The Main Extractor procedure 330 begins by initializing and starting a timer called the Message Timer (MsgTimer) (332). The Message Timer is typically programmed to expire in 1 or 2 seconds, although the timeout period is configurable to virtually any value. Next, the extractor procedure reads a record in the MAT (334). If the MAT record is an audit record for an RDF protected object, is a commit/abort record for any transaction, or is a "Stop Updaters" record, the audit record is modified by inserting the MAT position of the current audit record into the audit record and by appending to the audit record a timestamp (336). The appended timestamp is the timestamp of the last transaction to complete prior to generation of the audit record in the MAT 204. Every time the extractor procedure encounters a commit or abort audit record, it moves a copy of the timestamp in that record into its local timestamp register 297. The value in the local timestamp register 297 is the timestamp that is appended to audit records so as to generate an audit image record, also known as an image record.

If the message buffer currently in use has room for the resulting audit image record (338) it is moved into the message buffer (340). However, if the audit records being processed is a Stop Updaters audit record (337), the extractor sends a copy of the Stop Updaters audit image record (with the embedded MAT position) to Monitor Process 240 using a waited message send (339). Only after the Monitor Process 240 responds to the waited message does the extractor process resume its operations and move the Stop Updaters audit image record into the current message buffer. The Monitor Process is described in more detail below in the section of the document entitled "RDF support for DDL operations."

After the current audit image record is moved into the current message buffer, the Extractor procedure continues processing the next record in the MAT at step 334.

If the message buffer currently in use is full (338), the values stored in the message sequence number register 295 and the timestamp register 297 are inserted into the Message Buffer's header 280 (342). The extractor procedure then transmits the message buffer to the receiver process (344). After transmitting the message buffer, the Message Buffer Status array 294 is updated to indicate that the message buffer just transmitted is not available for use. In addition, the Message Timer is cleared and restarted, and the Message Sequence Number in register 295 is increased by one (346). Finally, the audit image record which did not fit in the last message buffer is moved into a next message buffer (348). If a next message buffer is not available, the extractor procedure waits until one becomes available and then moves the audit image record into it. Then the Extractor procedure continues processing the next record in the MAT at step 334.

When the audit record read (334) from the MAT 204 is not an audit record for an RDF protected table, is not a transaction commit/abort record and is not a "Stop Updaters" record, the audit record is ignored and the next audit record (if any) in the MAT is read (334).

The purpose of the Message Timer is to ensure that audit image records are transmitted to the receiver process in a timely fashion, even when the rate at which audit records are generated for RDF protected tiles is low. Referring to Figure 6D, when the Message Timer times out the Message Timer procedure 360 first checks to see if the current Message Buffer is empty (i.e., contains no audit image records) (362). If so, a current timestamp indicative of the current time is inserted into the Message Buffer header 280 (364). If not, the timestamp value from the last commit/abort record, stored in timestamp register 297, is inserted into the Message Buffer header (366). Then the current Message Sequence Number is inserted in the Message Buffer header (368) and the Message Buffer is transmitted to the receiver (370). After transmitting the message buffer, the Message Buffer Status Array 294 is updated to indicate that the message buffer just transmitted is not available for use, the Message Timer is cleared and restarted, and the Message Sequence Number in register 295 is increased by one (372).

When the extractor process receives a reply from the receiver process acknowledging receipt of a message buffer (374), if the reply message indicates the message buffer was received without error, the Message Buffer Status Array 294 is updated to indicate that the message buffer identified in the reply message is available for use (376).

If the reply message received by the extractor process from the receiver process has a message value indicating an error condition, the receiver process is requesting that the extractor and receiver resynchronize. The receiver process sends an error condition reply message whenever (A) a message with an out-of-sequence Message Sequence Number is received, and (B) whenever the receiver process starts up after a failover or return of control back to the primary receiver process from the backup receiver process (sometimes called a CheckSwitch). When the extractor process receives an error condition reply message from the receiver process, it waits for any pending message acknowledgment replies to be received for any other message buffers transmitted prior to receipt of the error condition reply message, and it ignores those reply messages (378). Then the extractor process calls the Extractor Volatile Initialization procedure (379) so as to resynchronize the extractor process with the receiver process.

Detailed Description of Receiver Process

The primary data structures used by the receiver process 232 in the preferred embodiment are shown in Figures 7A-7D. As stated earlier, the receiver process durably stores a context record 270 on a nonvolatile (disk) storage device 272 on a periodic basis (e.g., once every 5 to 25 seconds). As shown in Figure 7A the context record includes two location values per auxiliary image trail 238, three location values for the master image trail 236 (including a TST/MIT restart position 390), and a StopUpdatersCnt count value 391.

As explained earlier, every audit record shipped to the receiver process 232 has a MAT position value inserted in it by the extractor process. When the receiver process moves audit image records into the image trail buffer 274 for an image trail, the MAT position of the last such record is stored in the appropriate slot of the context record 270.

Furthermore, each image trail buffer 274 is written to the corresponding disk file only (A) when the image trail buffer 274 is full (i.e., contains 28K of data) or (B) when the receiver process performs a periodic flush operation. Each time data from any image trail buffer 274 is written to disk, the disk file location for the next write to the image trail file (i.e., the disk address for the current end of the image trail file) is stored in the appropriate slot of the context record 270.

The restart MIT position 390 in the context record is derived from the updater status table 400, as will be described below. Basically, the restart MIT position is a position in the Master Image Trail at which the receiver process 232 should start processing audit records for rebuilding the transaction status table 144.

The StopUpdatersCnt 391 is a count value that is incremented each time the receiver encounters and StopUpdaters record in a received message buffer whose MAT value is higher than the MAT position for at least one image trail.

The image trail buffer status array 392 stores one pointer value for each image trail buffer 274. That pointer value is the buffer location for the next audit record to be written to the image trail buffer, and is updated every time an audit

record is written to the image trail buffer to point to the next available position for an audit record in that buffer. Whenever the buffer pointer value reaches a 4K byte boundary, certain information (described below) is written into the next block header 394 in the image trail buffer 274.

The receiver process also stores a "Next Message Sequence Number" 396, a "restart MAT position" 398, and an "ExpectStopUpdate" flag 399. The Next Message Sequence Number 396 is the message sequence number the receiver expects to see in the next message buffer receiver and is normally incremented by one after each message buffer is received. The restart MAT position 398 is the lowest of the MAT position values stored in the context record 270. The ExpectStopUpdate flag 399 is a flag set in response to a special "Expect Stop Update" message from the Monitor process just prior to a StopUpdaters audit record being moved by the extractor process into its current message buffer.

Referring to Figure 7C, the updater status table 400 maintained by the receiver process stores the following information for each updater process: a value 402 that identifies the associated image trail, the StopPoint location 404 currently assigned to the updater, the last MIT position 406 reported by the updater process to the receiver process, the last image trail file position 408 reported by the updater process to the receiver process, and a status value 409 (e.g., open or closed). The meaning of these fields will be explained below.

Referring to Figure 7D, the transaction status table 144 maintained by the receiver process includes for each commit/abort record processed in the master image trail 236 a transaction ID 410, a status value 412 (i.e., commit or abort), and the MIT position 413 of the commit/abort record. The receiver also maintains a MIT position value 414 representing the last record in the MIT processed by the receiver process for purposes of updating the transaction status table 144.

Referring to Figures 8A-8G, and the pseudocode in Appendix 2, the receiver process 232 works as follows.

Referring to Figure 8A, the Receiver Startup Procedure 440 is called whenever the receiver process 232 or its backup is started, as in the case of a failover or a transfer of control back to the primary receiver process 232 from the backup receiver process. The Startup procedure begins by creating a backup process (442). The startup procedure then performs a "static initialization" of the receiver process (444), which means that all static data structures used by the receiver process are allocated and initialized. Then a checkpoint operation is performed in which a takeover location is transmitted to the backup receiver process (446). The takeover location is, in essence a program address, and in the preferred embodiment the takeover location is the program location at which execution of the Receiver volatile initialization procedure 450 begins. Finally, the Receiver Startup procedure calls (448) the Receiver Volatile Initialization procedure 450.

Referring to Figure 8B, the Receiver Volatile Initialization procedure 450 is called during startup by the Receiver Startup procedure 440. The Receiver Volatile Initialization procedure 450 begins by reading the last stored Receiver context record from disk and using it as the receiver's current context record in volatile memory (452). Then the Receiver Volatile Initialization procedure allocates and initializes all volatile data structures (454) used by the Receiver process, including the image trail buffers 274, the image trail buffer status array 392, the updater status array 400 and the transaction status table 144. Then the Receiver Volatile Initialization procedure sets the receiver's Expected Message Sequence Number to "1" (456). This will force the receiver and extractor to resynchronize, unless the extractor is starting up at the same time such as in response to a "Start RDF" command. Finally, the Volatile Initialization procedure calls (458) the main Receiver procedure 460.

Referring to Figures 8C-8F, the Main Receiver procedure 460 includes a subprocedure 470 for periodically saving the receiver's context record. This subprocedure is called every M seconds, where M is preferably a value between 5 and 25 and is typically set to 15. The first step (472) of the context save procedure durably stores the contents of each image trail buffer to its associated image trail file and resets all image trail buffer pointers (in IT buffer status table 392) to point to the beginning of each image trail buffer 274. The subprocedure also copies the current MIT end of file position from the context record into the first block header 394 in each of the image trail buffers. Then the subprocedure waits for the writes to the image trail files to complete (472).

Once the image trail buffer contents have been written to disk, the receiver's context record is durably stored on disk using a WaitUntilComplete write operation (474) (i.e., the step does not complete until the associated disk process sends a completion message back to the receiver process). Finally, the image trail file locations in the context record are copied into the StopPoint locations in the Updater Status table 400 (476). This last step allows the Updaters to process records up through the last audit image record stored on disk prior to the last context save operation.

The receiver process 232 is a "passive" process in that it does not initiate messages to other processes. Rather it only responds to messages from the extractor process 230, messages from the updater processes 234, and from the monitor process 240.

When a message is received from the extractor process (462), if the message is a Resynch request message, the receiver determines which of the MAT positions listed in the context record is lowest (464), and sends a Resynch Reply message to the extractor with the determined lowest MAT position embedded in the reply message (466).

If the received extractor message is a message buffer message, the message sequence number (denoted Mes-

sage.SequenceNumber) in the received message is compared with the locally stored Next Message Sequence Number (468). If the received message sequence number is not equal to the locally stored Next Message Sequence Number, the received message buffer is discarded (480) and an Error Reply message is sent to the extractor (482).

If the received message sequence number is in sequence, the locally stored Next Message Sequence Number is incremented by one (484) and a "Message Buffer OK" reply is sent to the extractor (484). A message buffer identifier is associated with the received message and is also associated with the reply message so that the extractor can properly update its message buffer status table by marking the acknowledged message buffer as available.

Next, all the audit records in the received message buffer are processed in sequence (490). First, the image trail associated with the record is determined (496). For data update audit records the associated image trail is determined by determining the database object updated on the primary system, determining the volume on which that object is stored, determining the updater responsible for replicating RDF protected files on that volume and then determining the image file associated with that updater. For commit/abort and StopUpdaters records, the associated image trail is the master image trail. Next, the MAT position (AuditRecord.MATptr) in the audit record is compared with the MAT position (IT.MATptr) for the identified image trail (498). If the audit record's MATptr is not larger than the image trail's MATptr, the audit record is ignored (500) because it has already been processed by the receiver.

Otherwise, if the record is a commit/abort audit record, it is moved into the MIT buffer 274-1 (492). If the "Stop Updaters" record has not been previously processed (as determined by step 498) (A) the StopUpdaters records is moved or copied into all the image trail buffers 274 record, and (B) the StopUpdatersCnt 391 count value in the receiver context record 270 is incremented (494).

Otherwise, the audit record is moved into the identified image trail buffer (502). However, if doing so would overflow a 4K byte block within the image trail buffer (504) special processing is required (see description of steps 510, 512 below). Furthermore, if moving the audit record into the image trail buffer would overflow the last block in the image trail buffer (506) the entire image trail buffer through the last 4K block is durably stored in the associated image trail file, the receiver's context record is updated to reflect the new end of file position for the image trail file, any data in the overflow block is moved up to the beginning of the image trail buffer, and then the buffer pointer for that buffer is reset to point to the next available slot near the beginning of the image trail buffer (508).

If either a 4K byte block has been overflowed, or the entire message buffer has overflowed, the current MIT file position (which is the MIT file position associated with the last audit record written to the MIT message buffer) is stored in the header for the next 4K block in the image trail buffer (510). Then the process of moving the current audit record into the image trail buffer is completed (512) and processing of the next audit record (if any) in the received message buffer begins at step 490.

If the received message buffer was empty (520), the receiver determines the highest of the MAT positions stored in the context record for all the image trails, which is equal to the MAT position of the last audit record received from the extractor in the last message buffer received that contained any audit records. Then an "RDF control record" is moved into all the image trail buffers (524). The RDF control record denotes (A) the determined highest MAT position, and (B) the timestamp value in the received message buffer's header.

If the received message buffer was not empty (520), but if one or more image trails received no audit records from the current message buffer (526), the receiver determines the highest of the MAT positions stored in the context record for all the other image trails (528), which is equal to the MAT position of the last audit record received from the extractor in the current message buffer. Then an "RDF control record" is moved into each image trail buffer that did not receive any audit records (530). The RDF control record denotes (A) the determined highest MAT position, and (B) the timestamp value in the received message buffer's header.

When a message is received from any updater process (540), the message is initially processed by extracting from the updater request message (A) the updater's identifier, (B) the updater's current image trail position, and (C) the last MIT position read by the updater from a block header in the updater's image trail (542). The extracted Image trail position and MIT position are stored in the Updater Status table (544). Using the updated information in the Updater Status table, the Restart MIT position in the receiver's context record is updated to be equal to the lowest MIT position for all the updaters (546). In addition, the receiver process deletes any image trail files eligible for deletion based on the updated image trail position values received from the updater (547).

If the received Updater request message is a StopPoint request message, the receiver reads the current StopPoint location for the updater and sends a reply message with that StopPoint location to the requesting updater (548).

If the received Updater request message is a Transaction Status request message, the receiver extracts from the request message the Transaction identifier for which the updater is requesting a status value (550). Then the transaction status table 144 is inspected to see if the status of that transaction is known (552). If not, the receiver waits until a commit/abort record for the specified transaction is received from the extractor (554). Once the transaction status of the identified transaction is known, a reply message is sent to the requesting extractor (556). The reply message includes (A) the current StopPoint for the updater, and (B) a set of transaction status records, starting with the record for the identified transaction. Preferably a substantial number of transaction records (e.g., 200 records, or all the status

records in the TST 144 after the identified record if less than 200) is sent with each reply so as to reduce the number of Transaction Status requests sent by updaters.

Detailed Explanation of Updater Processes

The primary data structures used by each updater process 234 in the preferred embodiment are shown in Figure 9. Each updater process durably stores a context record 570 on a nonvolatile (disk) storage device on a periodic basis (e.g., once every 2 to 10 minutes, with 5 minutes being preferred). As shown in Figure 9 the context record includes a saved image trail file position 571, a saved last MIT position 572, and a StopUpdaterCompleted flag 574. Each updater also stores in volatile memory

- a current image trail file position 575,
- a "Last Read MIT Position 576, which is the last MIT position read by the updater from a block header in the updater's image trail,
- a Highest Timestamp value (577), equal to the highest timestamp of any audit record processed by the updater,
- a StopPoint image trail file position (578),
- a scratch pad (579) for processing audit records, and
- a local transaction status table (580) provided to the updater by the receiver process.

The Highest Timestamp value 577 is not used by the RDF procedures discussed here. However, it is accessible by procedures executed by the Monitor process 240 for monitoring how far the updaters are running behind the TM/MP 202, and thus how long it would take the RDF system 220 to synchronize the backup database 124 with the primary database 108 if all transactions on the primary system were to stop.

Referring to Figures 10A-10E, and the pseudocode in Appendix 3, the updater processes 234 work as follows.

Referring to Figure 10A, the Updater Startup Procedure 600 is called whenever any updater process 234 is started. The Updater Startup procedure begins by creating a backup process (602). The startup procedure then performs a "static initialization" of the updater process (604), which means that all static data structures (such as a map of primary volumes to backup volumes) used by the updater process are allocated and initialized. Then a checkpoint operation is performed in which a takeover location is transmitted to the backup updater process (606). The takeover location is, in essence a program address, and in the preferred embodiment the takeover location is the program location at which execution of the Updater Volatile Initialization procedure 610 begins. Finally, the Updater Startup procedure calls (608) the Updater Volatile Initialization procedure 610.

Referring to Figure 10B, the Updater Volatile Initialization procedure 610 is called during startup by the Updater Startup procedure 600. The Updater Volatile Initialization procedure begins by reading the last stored Updater context record from disk and using it as the updater's current context record in volatile memory (612). Then the Updater Volatile Initialization procedure allocates and initializes all volatile data structures (614) used by the Updater process, including the scratchpad 579 and local transaction status table 580. Then the Updater Volatile Initialization sends a StopPoint request message to the receiver with the updater's current image trail position and MIT position embedded in the message, and stores the StopPoint value in the resulting reply message in its local StopPoint register 578. Finally, the Volatile Initialization procedure calls (618) the main Updater procedure 620.

Referring to Figures 10C-10E, the Main Updater procedure 620 includes a subprocedure 630 for periodically saving the updater's context record. This subprocedure is called every K minutes, where K is preferably a value between 2 and 10 and is typically set to 5. The first step (632) of the updater context save procedure 630 is to wait for all currently pending redo operations to complete so that the data stored in the updater's backup disk volume is known to reflect all updates through the last audit record read by the updater. Then the updater context save procedure copies the current Image Trail File Position 575 and the current Last Read MIT Position 576 into the corresponding fields of the context record 570, and durably stores the context record 570 on disk (634).

Referring to Figures 10D and 10E, the primary job of the Main Updater procedure 620 is to process audit image records in its image trail. At step 622 it reads the next audit record, if any, in the image trail. If, while reading that record, a 4K byte block header is encountered (624), the MIT position denoted in that header is moved into the updater's current MIT position register 576 (626).

If the audit record just read is an "RDF Control" record, the Highest Timestamp register 577 is updated by storing in it the timestamp in the RDF control record (628). Then processing resumes with next audit record (622).

If the audit record just read is a "Stop Updaters" record, the "StopUpdaterCompleted flag 574 in the Updater context record 570 is set to True (640) and the Updater context save procedure 620 is called (642). The StopUpdaterCompleted flag 574 is read by the Monitor process to ensure that all Updaters have stopped and that all have processed their image trails through the StopUpdaters record (as opposed to stopping due to a failure). Then the updater's backup process is terminated and the updater process itself terminates (644). The updater process will startup again after the

operator of the RDF system performs on the remote backup system the DDL operation which created the Stop Updaters audit record and then enters either the "Start Update" or "Takeover" command.

If the audit record just read is a database object update (i.e., data update) record, the Highest Timestamp register 577 is updated by storing in it the timestamp from the audit image record (646). The transaction identifier is extracted from the audit record (648). Then the local transaction status table 580 is inspected to see if it contains a status record for the identified transaction (582). If not, a Transaction Status Request message is sent to the receiver (652). When a reply message is received, the local transaction status table 580 is replaced with the transaction status table embedded in the reply message, and the StopPoint stored in StopPoint register 578 is replaced with the StopPoint value embedded in the reply message.

Once the local status transaction table contains a status record for the identified transaction, the transaction status of that transaction is determined (654). If the transaction status is "commit," a redo of the database object update operation noted in the audit record is initiated against the backup database file (656). If the transaction status is "abort," the audit record is not further processed. In either case, processing resumes with next audit record (622).

When the attempt to read a next audit record (622) encounters an audit record at or beyond the StopPoint value in StopPoint register 578, a StopPoint request message is sent to the receiver (660) to determine whether the StopPoint for the updater has been advanced. When a reply message is received, the StopPoint value in the received message is compared with the locally stored StopPoint value (622). If the two are equal, the updater 234 cannot process any further audit image records. As a result, the updater waits for W seconds (664), where W is preferably a value between 1 and 10 and is typically set to 5, and then sends another StopPoint request message to the receiver (660). This continues until a new StopPoint value is received from the receiver. At that point the locally stored StopPoint value in StopPoint register 578 is replaced with the StopPoint value in the received reply message, and then processing of audit image records resumes at step 622.

Online DDL Operations

For the purposes of this discussion, an object is defined to be a file, database table or other encapsulated computer resource accessed by a program as a unitary structure. In the context of the preferred embodiment, objects are database tables or indexes. In other implementations of the present invention, objects may be other encapsulated computer resources that the end user accesses indirectly through validated methods (i.e., programs) designed specifically to access those computer resources.

Figure 11A shows the data structure of a typical database table 720-1. The table 720-1 includes a file label 760, which is essentially a compactly stored copy of the catalog information for the database table, which represents the table's schema as well as other information not relevant here. Next, the table includes a primary key B-tree 762. The use of B-trees in database files is well known to those of ordinary skill in the art of database management systems. Next, the table has a data array 763 organized in rows and columns. The rows of the table are often called "records" 764.

In the context of embodiments of the present invention, every database table has a primary index. The value of the primary index for a particular record is called the primary key, and the primary key is typically equal to either (A) the value in one field (i.e., column), (B) the concatenation of the values in a plurality of columns, or (C) a computed function of the values in one or more columns. The set of columns used to generate the primary key are represented by a vector herein called the PrimaryKeyColumns vector. There is a "Create_PrimaryKey" function for every database table, represented as follows:

$$\text{PrimaryKey} = \text{Create_PrimaryKey}(\text{BaseTable}(\text{RecPtr}), \text{PrimaryKeyColumns})$$

where RecPtr is a pointer to a database table record.

It is often the case that an application program needs to access a database table in accordance with a set of column values, at least some of which are not included in the primary index. When that is the case, a Create Index procedure can be used to create an efficient alternate access path to the database table by ordering data according to the values in any specified set of columns. That ordering is represented by an "Alternate Index," which is typically implemented as a separate data structure from the associated database table.

Figure 11B shows the data structure of an alternate index 770. The alternate index 770 includes a file label 772, which includes a compactly stored copy of the catalog information for the index. The alternate index also includes an alternate key B-tree 774 and then a data array 776 organized in rows and columns. The data array has two sets of columns herein called the AltKey columns and the PrimaryKey columns.

The rows of the data array 776 are called records 778, and each row of the Alternate Index corresponds to one record of the associated database table. Furthermore, each row of the Alternate Index has two fields: one represents the alternate key value for the corresponding database table record, and one represents the Primary Key value for the same database table record.

The set of columns used to generate the alternate key for a particular database table are represented by a vector

herein called the AlternateKeyColumns vector. There is a "Create_AltKey" function for every alternate index of any database table, represented as follows:

AltKey = Create_AltKey(BaseTable(RecPtr), PrimaryKeyColumns)

where RecPtr is a pointer to a database table record.

For the purposes of the Move Partition, Split Partition and Move Partition Boundary operations, a database index may be viewed in the same way as a database table.

Figure 11C shows the data structure of an Audit Trail 738-1. The Audit Trail includes a file label 782 and a set of sequentially generated and stored audit entries 784. Most audit entries denote a database table record event, such as an addition, deletion or modification of a specified database table record in a specified database table. Other audit entries indicate that a specified transaction has committed or aborted. Yet other types of transactions, such as transactions involving restructurings of the database tables (e.g., adding an additional alternate index to a database table or changing the number of partitions for a database table), while still other audit entries may denote other events not relevant to this document.

Overview of Database Table Alteration Procedures

Figure 12 is a conceptual representation of the procedure for modifying a database table or index. The commands for making database table or index alterations are called data definition language (DDL) statements. In the preferred embodiment, the DDL statements used are not changed, except for the use of an "ONLINE" option in the DDL statements to indicate that the DDL operation is to be performed while minimizing the impact of the DDL operation on user transactions. When the ONLINE option is specified, the preferred embodiment of the present invention changes how the SQL catalog manager executes the Move Partition, Split Partition, Move Partition Boundary and Create Index commands.

The procedure for modifying a database table's structure begins with a user or operator entering a DDL statement, specifying an alteration in the schema of a specified object 790 which is typically either a database table, an index, or a partition of a database table or index. The specified object is accessed to read its file label, and if the command requires generating a new object 792, the new object is created. If the command involves movement of data between the first object 790 and a second object 792 that already exists, the file label of the second object is also accessed.

In the first phase of execution of the "online DDL command," a "dirty copy" of data from the first object 790 into the second object 792 is made while user transactions against the first object (and against the second object if it existed prior to the command) are allowed to continue unimpeded by the execution of the DDL command. In particular, a "dirty copy" is made by accessing all the data records of the first object that are the subject of the DDL command, using "browse" access, and generating corresponding records in the second object 792 until the last of the relevant data records in the first object 790 have been accessed.

Browse access is a form of read access that allows database records to be read through any record locks that may be in place. Thus, browse access will sometimes result in the reading of records which are in the midst of being modified. Also, user transactions against the first object may delete or modify records in the first object while or after they are accessed by the first phase of the DDL command execution. As a result, the records created in the second object 792 may require a certain amount of correction.

In the second phase of execution of the online DDL command, a partial clean-up of the records in the second object is performed by accessing records in the audit trail 794 associated with the first object 790 and performing corresponding "redo" operations against the second object. During this second phase user transactions against the first object (and against the second object if it existed prior to the command) are allowed to continue unimpeded by the execution of the DDL command.

More specifically, before the dirty copy in the first phase of the online DDL command is executed, a pointer called AuditPtr is set to the end of the audit trail 794. During the second phase, each audit record beginning with the one referenced by the AuditPtr is inspected. If the audit record is relevant to operations performed on the first object 790, or relevant to the subset of records of the first object that are the subject of the DDL command, then a redo operation is performed against the second object 792 based on the information in the audit record. In the case of a Create Index command, the redo operation uses new audit records that are generated based on the audit entries found in the audit trail because the second object has a very different structure than the first object; in the case of the other three online DDL commands, the redo operation uses the audit record as found in the audit trail except that the audit record is modified prior to the redo operation so as to reference the second object instead of the first object, and the redo operation is performed against the second object.

The second phase continues until the end of the audit trail is reached.

In the third phase, a lock is requested against the first object and second object (if it existed prior to the DDL command), and when that lock (or those locks) is (are) granted all user transactions other than browse access trans-

actions are blocked until the third phase is completed. During the third phase, any new audit records in the audit trail 794 associated with the first object 790 are accessed and corresponding "redo" operations are performed against the second object 792. Next, if the first object is a portion of a larger database table or index, a lock against the entire associated database table associated with the first object is obtained while catalog and file labels are updated. In particular, the catalog entries and the file labels associated with the first and second objects are updated to reflect the results of the DDL operation. Furthermore, the catalog entries and file labels of all objects which incorporate schema information made invalid by the DDL operation are also updated. Then the lock on the first and second objects (and the lock, if any, on the associated database table or index) is released, enabling user transactions against the first and second objects to resume.

Finally, if the DLL command requires deletion of the first object or deletion of a range of records in the first object, that deletion operation is performed in a way that permits concurrent transaction activity.

It is noted that the implementation of database reconfiguration operations using well established, optimized, database recovery process procedures makes the database reconfiguration operations efficient in terms of computation time and in terms of the computer resources used.

RDF Support for Online DDL Procedures

Referring to Figure 3, a special audit entry, herein called a "Stop Updaters" audit record, is made when any online DDL operation successfully completes. This audit entry is stored in the master audit trail 204. As described above, when the extractor process 230 reads a Stop Updaters record in the MAT 204, it stops processing audit records, and sends a Stop Updaters message to the monitor process 240 indicating (A) that a Stop Updaters audit entry was encountered and (B) the MAT position of that audit entry. The Stop Updaters message is a waited message, causing the extractor process to be suspended until it receives a reply from the monitor process.

Referring to Figure 13, the monitor process has a context record 800 that includes two fields: an OnlineDDLcnt field 802 and an associated MAT value 804. The OnlineDDLcnt field 802 is equal to the number of online DDL operations that have been performed on the primary system that have yet to be performed on the remote backup system. The MAT value 804 is the MAT position of the Stop Updaters audit record with the highest MAT position.

Referring to Figure 14, when the Monitor process 240 receives a Stop Updaters message from the extractor, it executes a StopUpdaters Message procedure 810. That procedure first compares the MAT position value in the message with the MAT position value in the Monitor's context record (812). If the message MAT value is larger, that indicates the Monitor has not previously been informed of the current Stop Updaters audit record, and the Monitor increments its OnlineDDLcnt counter in its context record (814). Then it stores the message's MAT position value in its context record (816) and durably saves its context record 800 on disk (818). Next, it sends a waited message to the receiver process indicating that a Stop Update will occur in the near future (820). When the receiver replies to that message, the Monitor sends a reply back to the extractor (822), thereby completing the waited message sent by the extractor.

If the message MAT position value is not larger than the MAT position value in the Monitor's context record (812), but the Monitor's context record has a nonzero OnlineDDLcnt value (830), then it is possible that a failure in either the extractor or receiver process has caused the receiver to not be aware of the upcoming Stop Update. The monitor therefore sends a waited message to the receiver process indicating that a Stop Update will occur in the near future (820). When the receiver replies to that message, the Monitor sends a reply back to the extractor (822), thereby completing the waited message sent by the extractor.

Referring to Figure 15, the Start Update command is the command used to start all the updaters. The Start Update procedure 850 used to execute the Start Update has been modified to support online DDL operations. When the Start Update command is issued, the Start Update procedure 850 (executed in the Monitor process 240) checks to see if all updaters are stopped (852). If not, an error message is sent back to the operator, since the updaters can be started only when they are all stopped. If the updaters are all stopped, the monitor checks to see if its OnlineDDLcnt is nonzero (854). If OnlineDDLcnt is equal to zero, then all the updaters are started (856).

If OnlineDDLcnt is greater than zero, the monitor next checks the context records of all the updaters to see if the StopUpdaterComplete flag is set in all the updater context records (858), which they will be if there have been no failures. When all the updaters have StopUpdaterComplete flags set, that means that all the updaters stopped execution in response to reading a StopUpdaters audit record. In that case, the monitor decrements by 1 the StopUpdatersCnt in the receiver's context record and decrements by 1 the OnlineDDLcnt in its own context record (860) and durably saves those context records on disk. Furthermore, the monitor turns off the StopUpdaterComplete flag in the context records of all the updaters (862) and then it starts all the updaters (856).

If OnlineDDLcnt is greater than zero, but the StopUpdaterComplete flag is not set in one or more of the updaters, that means that those updaters stopped prior to reading a StopUpdaters audit record. It is essential that all the updaters process all audit records up to the StopUpdaters audit record before the DDL operation performed on the primary system is repeated on the backup system. Therefore, when this condition is detected, the monitor will start just the

updaters whose StopUpdaterComplete flags are not set (870). Those updaters should then process audit records until they reach the StopUpdaters audit record that caused the other updaters to stop. At that point, the operator will execute the DDL operation against the remote backup database, and then will re-enter the Start Update command.

Referring to Figure 16A, the receiver procedure for responding to an Expect Stop Updaters message 880 works as follows. When an Expect Stop Update message is received from the monitor process, the receiver sets its Expect-StopUpdate flag in volatile memory to True (882) and then sends a reply back to the monitor process so as to complete the waited message initiated by the monitor process.

Referring to Figure 16B is a flowchart of the receiver procedure 890 for responding to an Updater Close message from an updater process. When an updater performs a normal termination, it sends a "Close message" to the receiver and monitor processes. The receiver responds to each close message by setting the status 409 of the sending updater in its updater status table 400 to "closed" (892). Then, if the ExpectStopUpdate flag is set to True (894) and all the updaters are closed (896), the receiver deletes its transaction status table and updater status table and resets the ExpectStopUpdate flag to False (896). In other words, only when an online DDL operation has been performed on the primary database does the receiver automatically delete its transaction status table and updater status table when all the updaters have closed. Otherwise, the receiver deletes its transaction status table and updater status table only in response to a Stop Update or Stop RDF command.

Alternate Embodiments

While the present invention has been described with reference to a few specific embodiments, the description is illustrative of the invention and is not to be construed as limiting the invention. Various modifications may occur to those skilled in the art without departing from the true spirit and scope of the invention as defined by the appended claims.

APPENDIX 1

PSEUDOCODE FOR EXTRACTOR PROCESS

5 Procedure: Startup Extractor()

Upon Startup of Extractor:

```

10 {
    Create a backup extractor process
    Perform static initialization of Extractor process
    Transmit a takeover location to the backup extractor process
    Call Volatile Initialization of Extractor()
15 }

```

Procedure: Volatile Initialization of Extractor()

```

20 {
    Initialize all volatile data structures in Extractor process
    Send request message to Receiver process asking for safe MAT to use
    Wait for Resynch Reply Message from Receiver process
25 Set MATptr to MAT value received from RDF Receiver process
    Set Extractor's Message Sequence Number to 1
    Call Extractor()
30 }

```

Procedure: Extractor()

Initialize and start Msg Timer /* times out once J every seconds (typically,
35 once every one or two seconds */

/* MATptr points to the next audit record to be processed */

Do Forever

```

40 {
    Read next audit record, if any
    If record is associated with an audited database table on an RDF
    protected volume, is a commit/abort record, or is a StopUpdaters
45 record
    {
        Insert MATptr into audit record
        If record is a StopUpdaters record
50 {
            Send StopUpdaters record to Monitor Process using a waited
            message send.
55 /* Extractor process waits for reply from Monitor */

```

}

Append to audit record a Timestamp from the last read commit or
abort record

If there sufficient room in Message Buffer to store audit record

Move augmented audit record into Message Buffer

Else

{

Insert Message Sequence Number in message buffer header

Insert timestamp from the last read commit or abort record in
Message Buffer header

Transmit Message Buffer to RDF Receiver

Update message buffer status array

Clear and restart MsgTimer

Increment Message Sequence Number

Move augmented audit record into next Message Buffer when it
becomes available /* either buffer is already available or
Extractor will wait until one becomes available via
acknowledgement from Receiver */

}

}

}

In Background:

Do Forever:

{

If Msg Timer has timed out

{

If Message Buffer contains no audit records

Insert current timestamp in Message Buffer header

Else

Insert timestamp from the last read commit or abort record in
Message Buffer header

Insert Message Sequence Number in Message Buffer header

Transmit Message Buffer to Receiver process

Update message buffer status array

Increment Message Sequence Number

Restart Msg Timer

}

Whenever a reply message is received from the Receiver process

```

{
  If message is a "Audit Message" acknowledgement
5    {
      If message is a "Message Buffer OK" reply message
          {
10              Update Message Buffer Status array to indicate that the
                  corresponding message buffer is available for use
          }
      Elseif message is an "Error / Need Resynch" reply message
15          {
              /* Need to resynchronize Extractor and Receiver Processes
*/
20              Wait for any other pending message acknowledgement
                  replies to be received (and ignore them)
                  Call Volatile Initialization of Extractor()
25              }
          }
      } /* End of reply message processing */
} /* End of background Do Loop */
30
35
40
45
50
55

```

APPENDIX 2

PSEUDOCODE FOR RECEIVER PROCESS

5 Procedure: Startup Receiver Process()

```

    {
      Create a Backup Receiver Process
10     Perform static initialization of Receiver process
      Transmit a takeover location to the Backup Receiver process
      Call Receiver Volatile Initialization()
15     }
  
```

Procedure: Volatile Initialization of Receiver()

```

    {
20     Recover last stored context record and use as current context record (in
        memory)
      Initialize all volatile data structures in Receiver process, including
25     transaction status table and image trail buffers
      Set Next Message Sequence Number to 1
        /* Setting SeqNo to 1 forces resynchronization */
      Call Receiver()
30     }
  
```

Procedure: Receiver()

```

35     {
      /* Save Receiver Context Subprocedure */
      Every M seconds /* M is typically between 5 and 25, with 15 being preferred
      */
40     {
        Store all IT buffers in associated Image Trail files
        Reset IT buffer pointers
45     Put current MIT file position in header of the first 4K block of each IT
        buffer
        Wait for writes to Image Trail files to complete

50     Durably store Context record using a WaitUntilComplete write operation
        Copy IT File locations from Context record to StopPoint locations in
        Updater Status Table
55     }
  
```

When an Extractor Resynch request message is received

```

{
  Determine lowest MAT position for all the image trails
  Send Resynch Reply Message with lowest MAT position to Extractor
}

```

Upon Receiving a Message Buffer

```

{
  If Message.SequenceNumber ≠ Next Message Sequence Number
  {
    Discard received message buffer
    Send "Error / Need Resynch" Reply Message to Extractor
  }
  Else
  {
    Increment Next Message Sequence Number
    Send "Acknowledgement" Reply Message to Extractor process
    For each data update record
    {
      Determine associated Image Trail
      If MAT value in data update record is greater than the MAT
        value for the last audit record moved into the associated
        Image Trail Buffer
      {
        Move data update record into the associated Image Trail
        buffer
      }
      Else
      {
        Ignore data update record /* it was already written to the
        image trail */
      }
    }
    Move all commit/abort records to buffer for Master Image Trail
    Copy each "RDF Stop Updaters" record into buffers for all image
    trails
  }
}

```

When any 4K block in the image trail buffer is full

Put current MIT position in header for the next 4K block

When any image trail buffer is too full to accept a record or new
block header

```

5      {
      Using an Unwaited Write, durably store the data in the buffer in
      the associated image trail file
      Reset Pointers to beginning of buffer
10     Resume updating block headers and moving audit record(s) into
      the image trail buffer
      }

```

15 Whenever a received message buffer is empty

```

      {
      Determine highest MAT of the "last MAT positions" for all the
20     Image Trails
      Write an RDF control record into all Image Trail buffers denoting
      the determined MAT position and the timestamp value in the
      message buffer's header
25     }

```

30 When receiver message buffer was not empty and any Image Trail
(an inactive Image Trail) receives no data update records from
the last received Message Buffer

```

      {
      Determine highest MAT of the "last MAT positions" for the other
35     Image Trails
      Move an RDF control record into the inactive Image Trail's buffer
      denoting the determined MAT position and the timestamp
40     from the Message Buffer header
      }

```

```

45     }
    } /* End of processing for a received message buffer of audit image
records

```

50 In Background:

```

      {
      Read all commit/abort records written to Master Image Trail and update
      Transaction Status Table accordingly
55     }

```

When any request is received from an Updater

```

{
5   Extract from Updater message the identity of the sending Updater, the
      Updater's current Image Trail position and MIT position
      Store Image Trail Position and MIT position in Receiver's Updater Status
      Table
10  Update Restart MIT position in Context Record (in memory)
      Delete any image trail files eligible for deletion based on updated image
      trail position values.

```

When request is a StopPoint request

/* StopPoint Request message is sent by an Updater when it restarts,
reaches its StopPoint, or reaches end of its audit trail */

```

{
20   Send reply message to Updater specifying its current StopPoint
      location
25   }

```

When the request is a Transaction Status request

```

{
30   Extract Transaction ID from Request
      Attempt to locate the status record in the transaction status table
      (TST) for that Transaction ID
      If record is not found
35   {
      Wait until commit/abort record is received from Extractor for the
      specified transaction
40   Update transaction status table with receiver commit/abort
      information
      }
45   Send Reply Message to requesting Updater with StopPoint position
      for Updater and N status records, starting at the located record
      (Send all status records to end of status table if number of
      records in table after the located record is less than N)
50   /* N is typically 200 or larger */

```

```

      }
    } /* end of processing of Updater Request */
55 } /* end of Receiver procedure */

```


APPENDIX 3

PSEUDOCODE FOR EACH UPDATER PROCESS

```

5  Procedure: Startup Updater()
    /* startup caused by either "Start RDF" or "Start Update" commands */
    {
10  Create a Backup Updater Process
    Perform static initialization of Updater process
    Transmit a takeover location to the Backup Updater process
    Call Volatile Initialization of Updater()
15  }

    Procedure: Volatile Initialization of Updater()
20  {
    Read last stored context record and use as current context record (in
        memory) /* includes current image trail position and MIT position */
    Initialize all volatile data structures in Updater process, including file and
25  buffer structures

    Send StopPoint Request Message to Receiver process, including
30  Updater's current Image Trail position and MIT position in the
        message
    Wait for reply
    Extract StopPoint position from Reply Message
35  Call Updater()
    }

40  Procedure: Updater()
    Do Forever
    {
    /* Save Updater Context Subprocedure */
45  Every K minutes /* K is typically between 2 and 10, with 5 being
        preferred */
    {
50  Wait for all pending redos to complete
        Durably store Context record using a WaitUntilComplete write
operation
55  }
    }

```

Read next record in Image Trail, but do not go beyond the last StopPoint
specified by the Receiver process

If at StopPoint

{
Do Until current Image Trail position = StopPoint

{
Once every W seconds, Send a StopPoint Request message to
the Receiver /* W is typically a value between 1 and 10 */

Wait for Reply

Extract StopPoint from reply message

}

}

If block header was encountered

Get MIT position from block header and store MIT position in
Updater's context record

BEGIN CASE (audit record type)

{

Case: record type = RDF Stop Updaters record

{

Perform normal context save operation, including durably storing
context record

Terminate backup Updater Process

Terminate Updater Process

}

Case: record = RDF control record

{

Update local HighestTimestamp with timestamp from record

}

Case: record = data update record

{

Update local HighestTimestamp with timestamp from record

Determine transaction ID of update

Check local transaction status table for status of identified
transaction

If transaction ID is not found

{

Send a Transaction Status Request to the Receiver process,
denoting the transaction ID not found in the local status
table. Include in the request message the current MIT

```

    position in the Updater's context record and the position
    of the current record in the Image Trail file
5    Wait until Receiver sends Reply Message with status data
    Replace local status table with status data received from
    Receiver (typically includes status records for 200 or
10   more transactions)
    }
    TxStatus = commit/abort status indicated in status table for
    identified transaction
15   If TxStatus=commit
    {
        Initiate a redo of the operation noted in the audit record
        against the backup database file. (The redo operation is
20   performed by disk process for the volume on which the
        file is located on the remote backup system.)
    }
25   } /* processing of data update record ends here */
    } /* End of Case Statement */
    } /* End of Do Loop */
30

```

Claims

35 1. A distributed computer database system, comprising:

a local computer system having a local database stored on local memory media, application programs that modify the local database, and a transaction manager that stores audit records in a local audit trail reflecting those application program modifications to the local database; each audit record having an associated audit trail position in said local audit trail;

a remote computer system, remotely located from the local computer system, said remote computer system having a backup database stored on memory media associated with the remote computer system;

a communication channel for sending messages between said local computer system and said remote computer system; and

45 a remote data duplication facility, partially located in said local computer system and partially located in said remote computer, for maintaining virtual synchronization of said backup database with said local database, including:

50 an extractor process executed by said local computer system that extracts audit records from said local audit trail and transmits said extracted audit records to said remote computer system;

a receiver process executed by said remote computer system that receives said audit records transmitted by said extractor process and distributes said audit records to a plurality of image trail files in said remote computer system; and

55 a plurality of updater processes executed by said remote computer system, wherein each updater process that reads said audit records in an assigned one of said image trail files and initiates redo operations of database modifications denoted in at least a subset of said read audit records against said backup database;

said local computer system including a manager for performing database restructurings while said application programs continue to modify said database;
 said transaction manager storing a Stop Updaters audit record in said local audit trail when each said database restructuring successfully completes;
 5 said extractor process including instructions for transmitting said Stop Updaters audit record to said remote computer system;
 said receiver process including instructions for moving a copy each received Stop Updaters audit record into all of said image trail files; and
 10 said updater processes each including instructions for stopping execution when said Stop Updaters audit record is read in said assigned one of said image trail files.

2. A distributed computer database system, comprising:

15 a local computer system having a local database stored on local memory media, application programs that modify the local database, and a transaction manager that stores audit records in a local audit trail reflecting those application program modifications to the local database; each audit record having an associated audit trail position in said local audit trail;
 a remote computer system, remotely located from the local computer system, said remote computer system having a backup database stored on memory media associated with the remote computer system;
 20 a communication channel for sending messages between said local computer system and said remote computer system; and
 a remote data duplication facility, partially located in said local computer system and partially located in said remote computer, for maintaining virtual synchronization of said backup database with said local database, including:

25 an extractor process executed by said local computer system that extracts audit records from said local audit trail and transmits said extracted audit records to said remote computer system; said extractor process including a plurality of message buffers and instructions for buffering groups of said extracted audit records together in said message buffers and transmitting said message buffers to said remote computer system, each transmitted message buffer having an associated sequence number, wherein sequentially transmitted messages buffers have associated sequence numbers that follow a predefined sequence;
 30 said extractor process instructions including instructions for continuing, after transmitting a first one of said message buffers to said remote computer system, to buffer groups of said extracted audit records in other ones said message buffers and for transmitting said other message buffers to said remote computer system;
 35 a receiver process executed by said remote computer system that

stores an expected message sequence number and a context record denoting a restart audit trail position value,
 40 receives said message buffers transmitted by said extractor process,
 responds to each received message buffer whose associated message sequence number does not match said expected message sequence number by transmitting an error message to said extractor process; and
 responds to each received message buffer whose associated message sequence number matches
 45 said expected message sequence number by sending a reply message to said extractor process to acknowledge receipt of each said message buffer, updating said expected message sequence number in accordance with said predefined sequence, distributing said audit records in said received message buffer to a plurality of image trail files in said remote computer system, and updating said restart audit trail position value based on said audit trail positions associated with said audit records
 50 in said received message buffer; and

a plurality of updater processes executed by said remote computer system, wherein each updater process that reads said audit records in an assigned one of said image trail files and initiates redo operations of database modifications denoted in at least a subset of said read audit records against said backup data-
 55 base;
 said extractor process including reply message instructions for responding to said reply message acknowledging receipt of one of said message buffers by said receiver process by enabling reuse of said one message buffer, and error reply message instructions for responding to said error message transmitted

EP 0 724 223 A1

by said receiver process by determining said restart audit trail position value stored by said receiver process and then extracting audit records from said local audit trail starting at said restart audit trail position value.

5

10

15

20

25

30

35

40

45

50

55

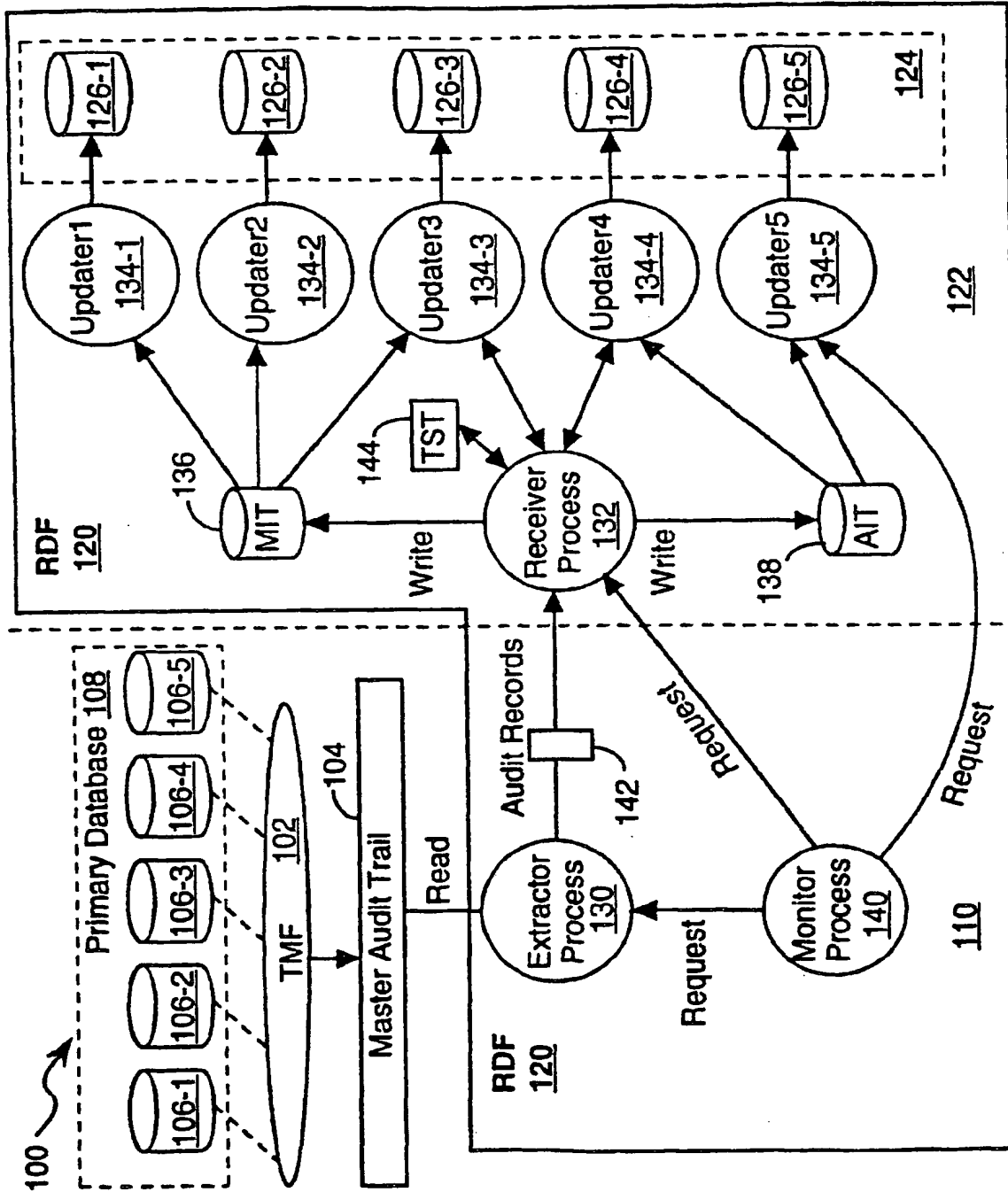


FIGURE 1

Prior Art

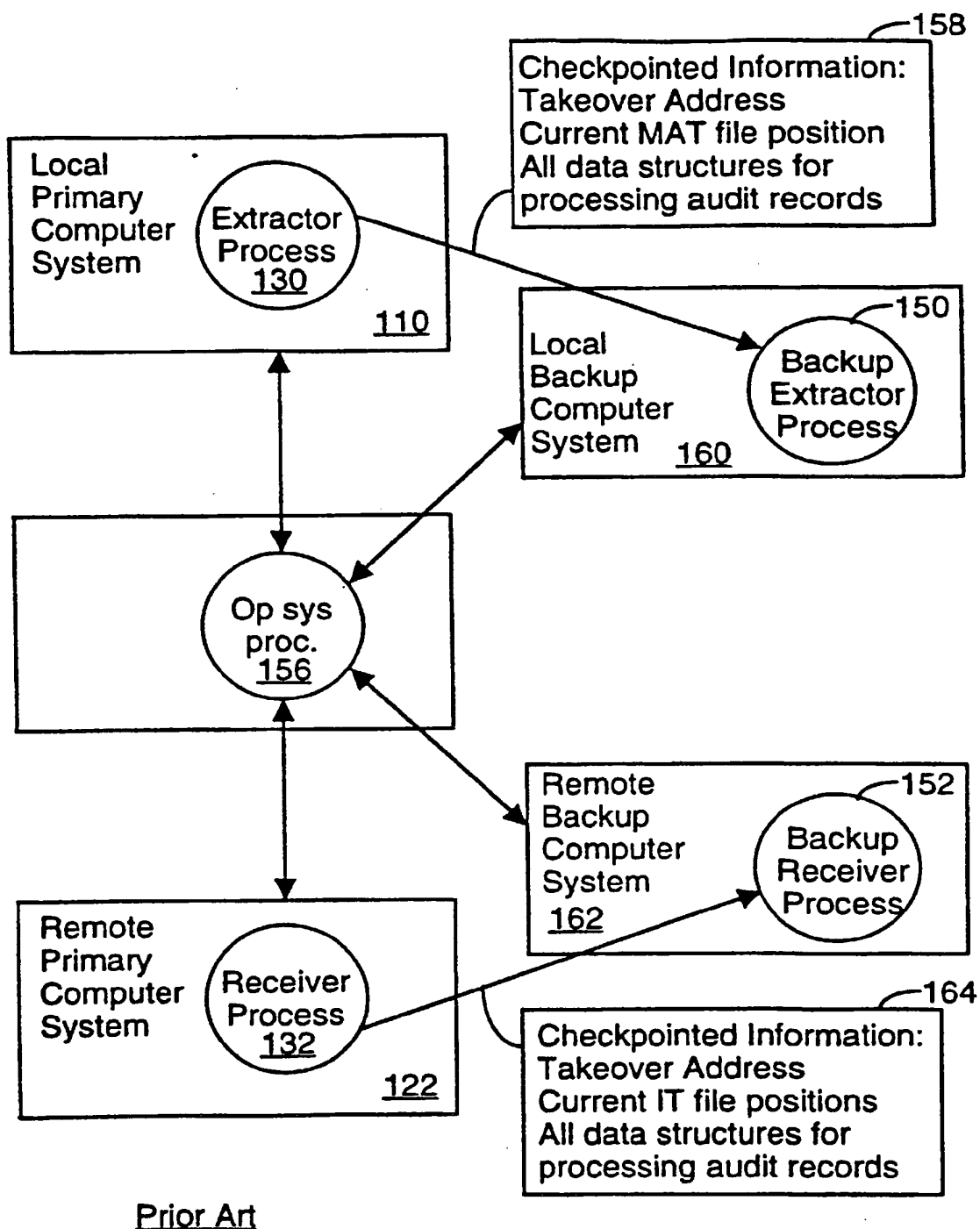


FIGURE 2

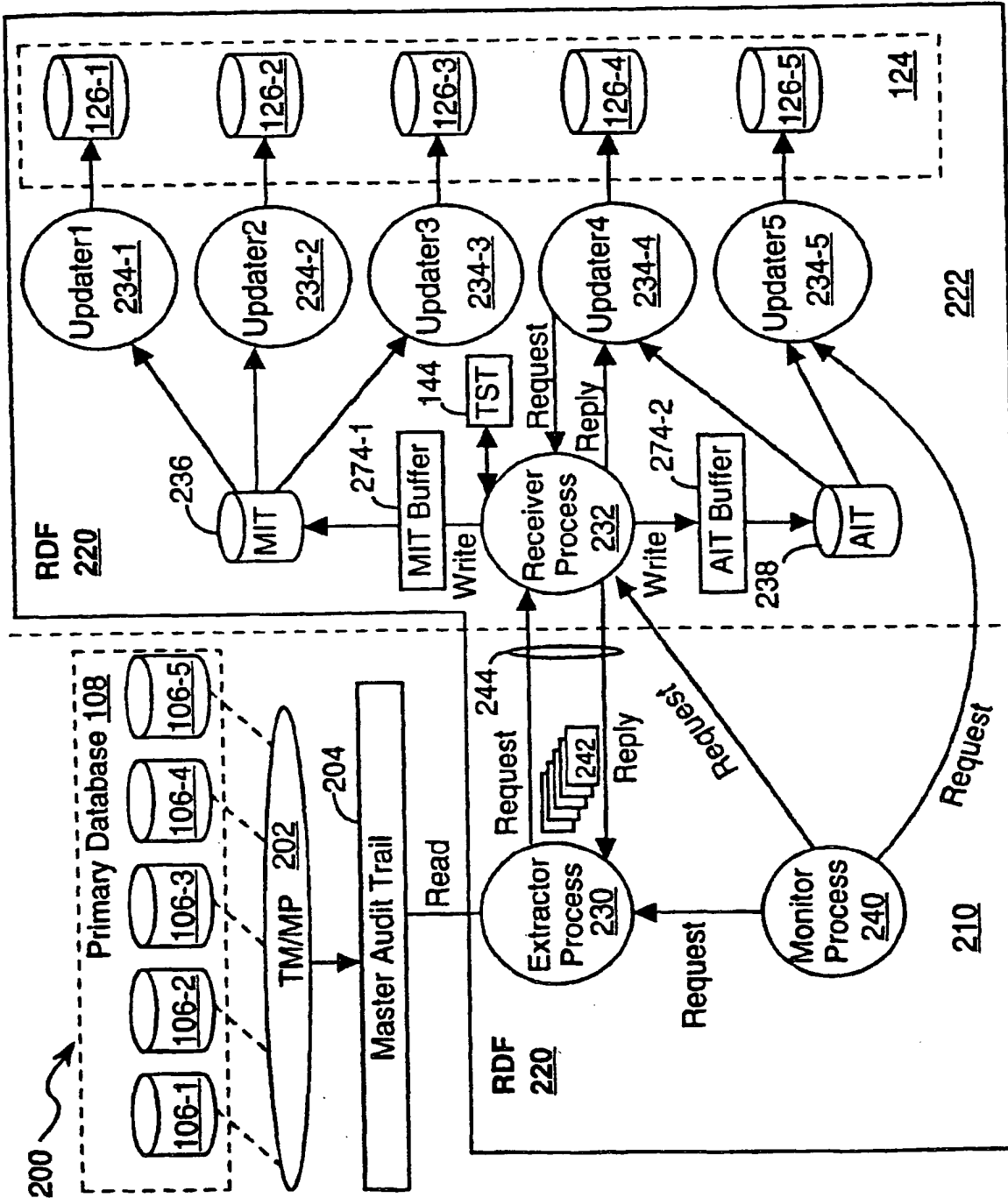


FIGURE 3

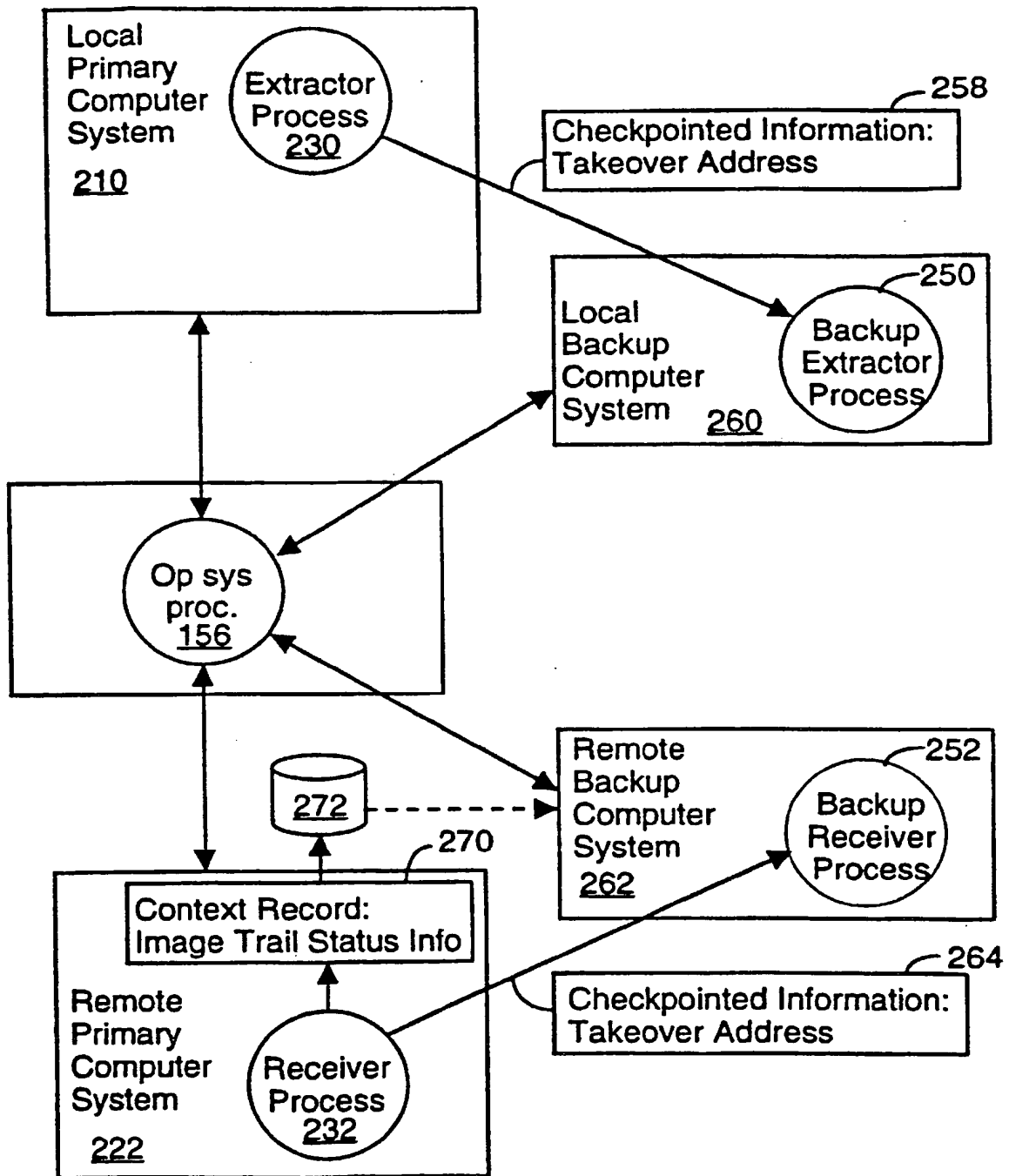


FIGURE 4

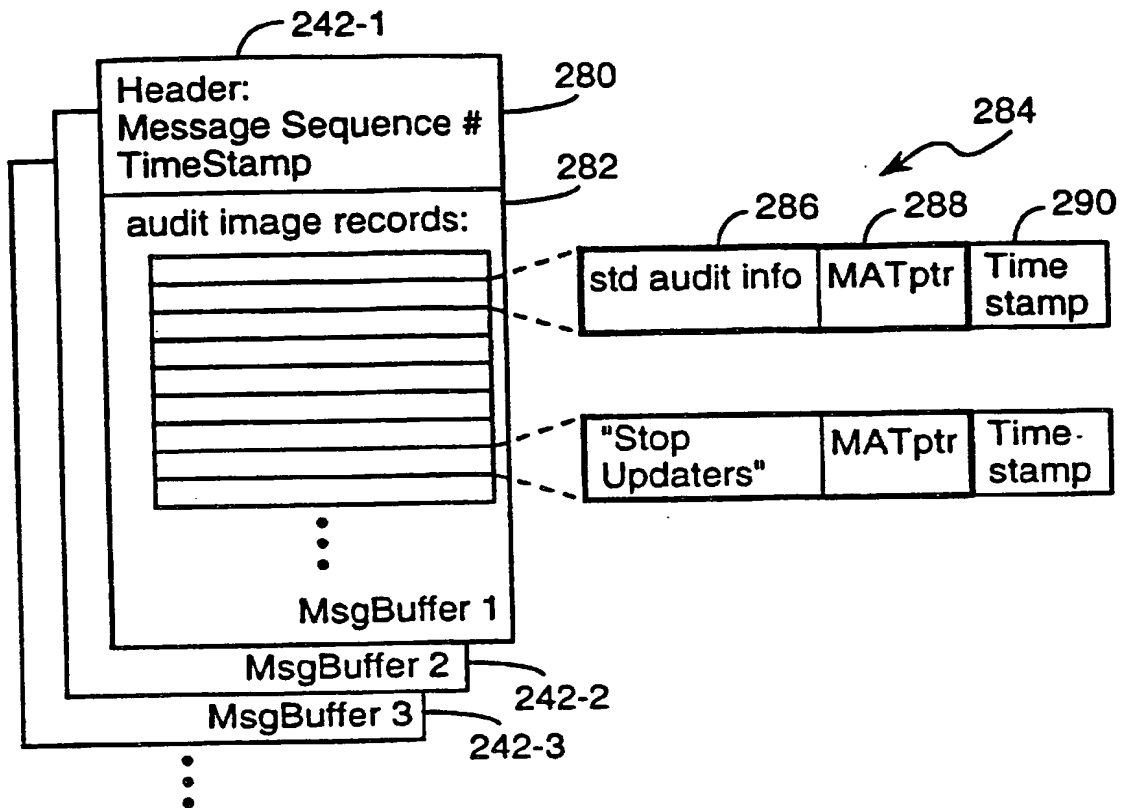


FIGURE 5A

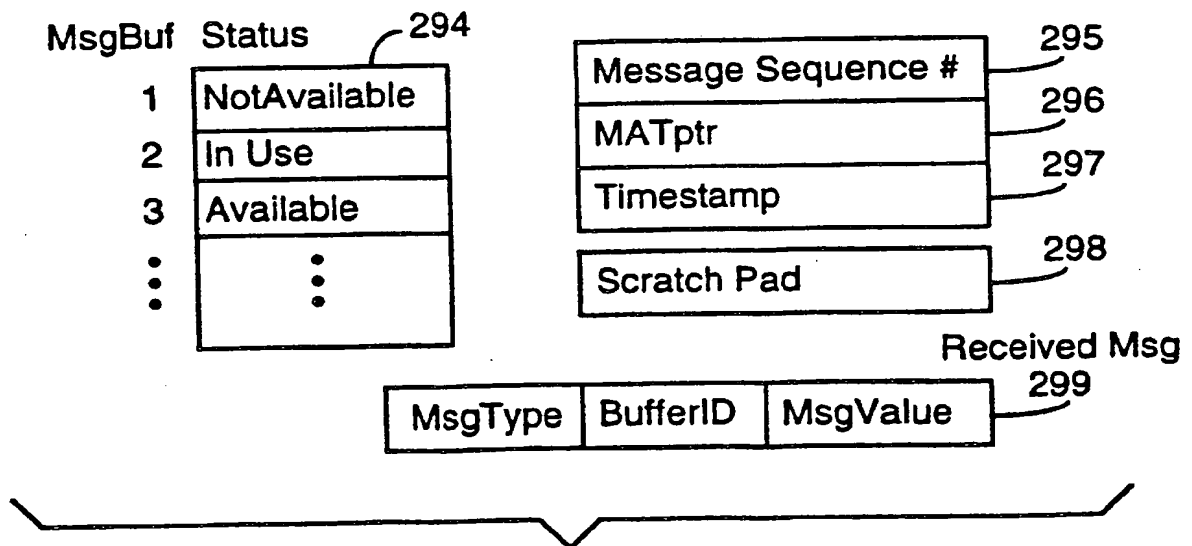


FIGURE 5B

Extractor Startup Procedure

300

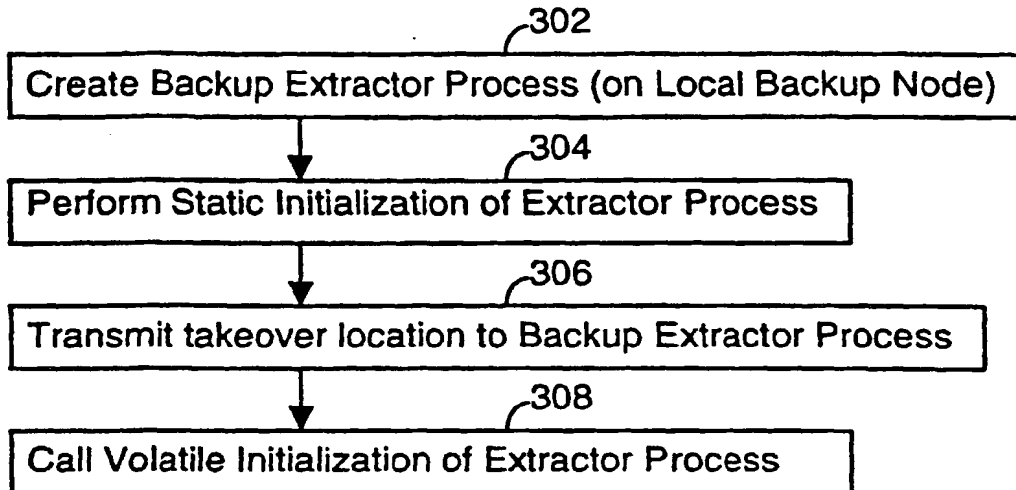


FIGURE 6A

Extractor Volatile Initialization Procedure

310

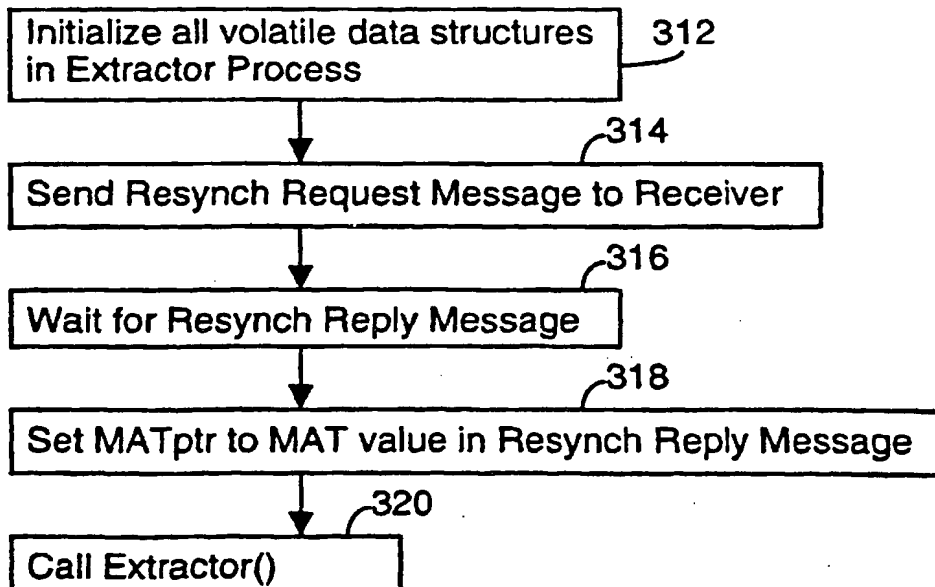


FIGURE 6B

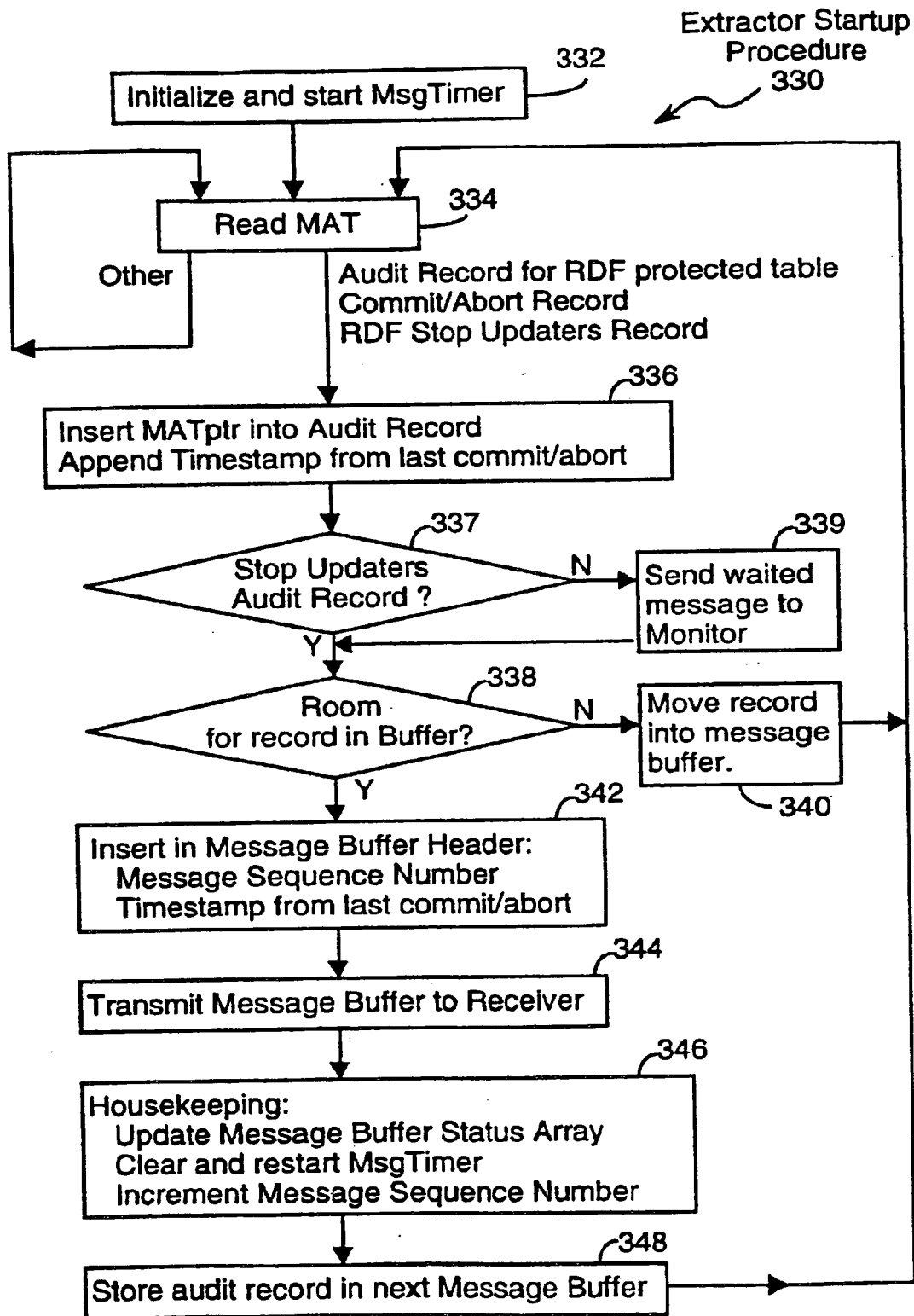


FIGURE 6C

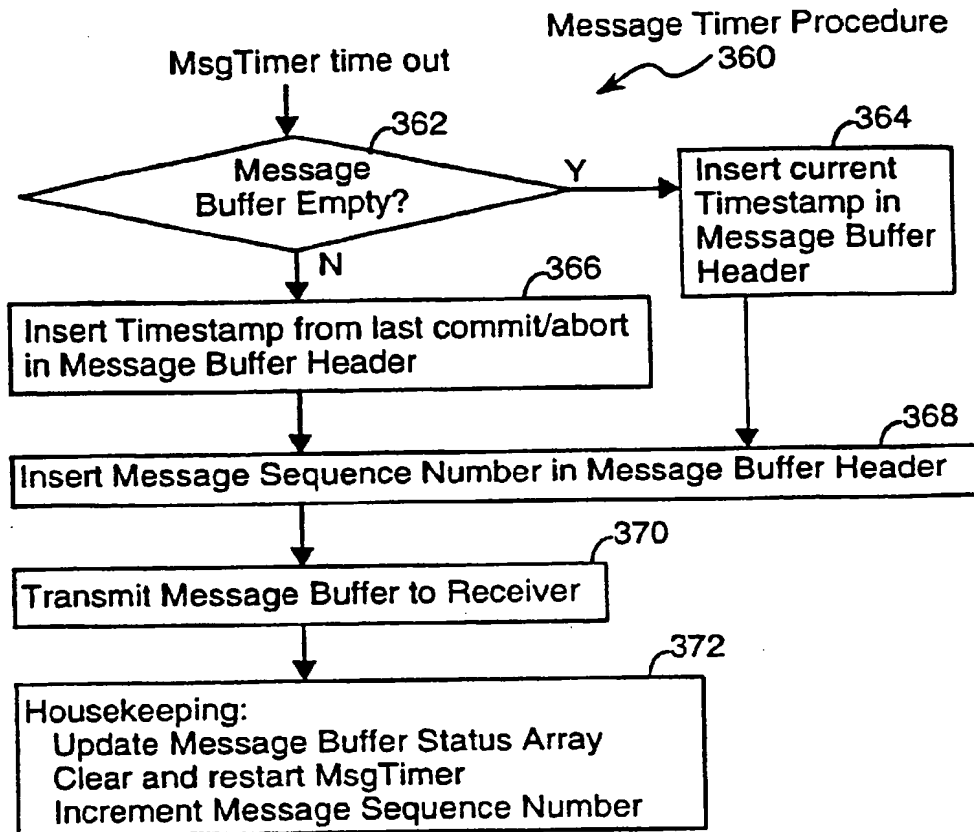


FIGURE 6D

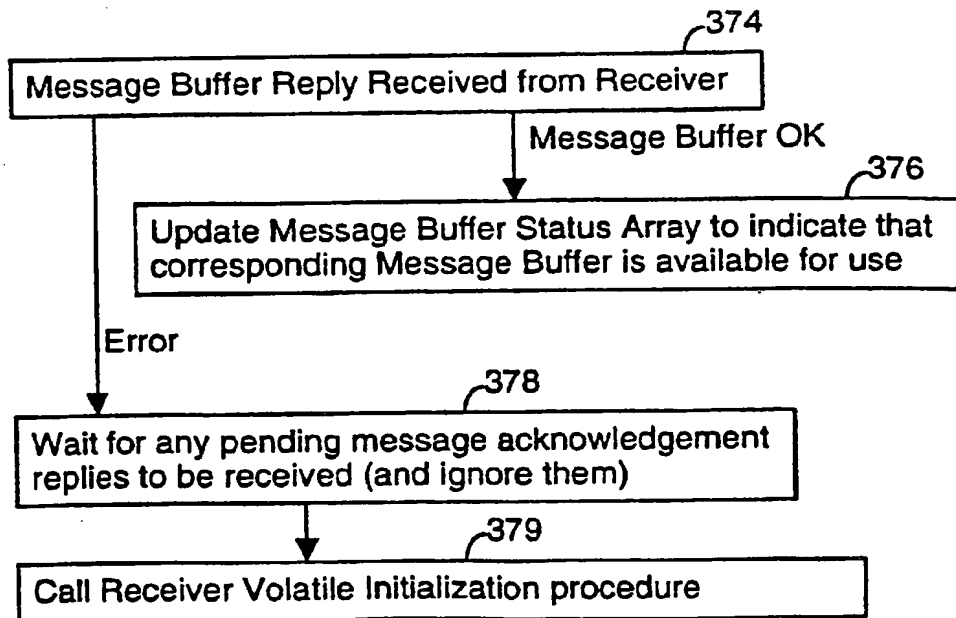


FIGURE 6E

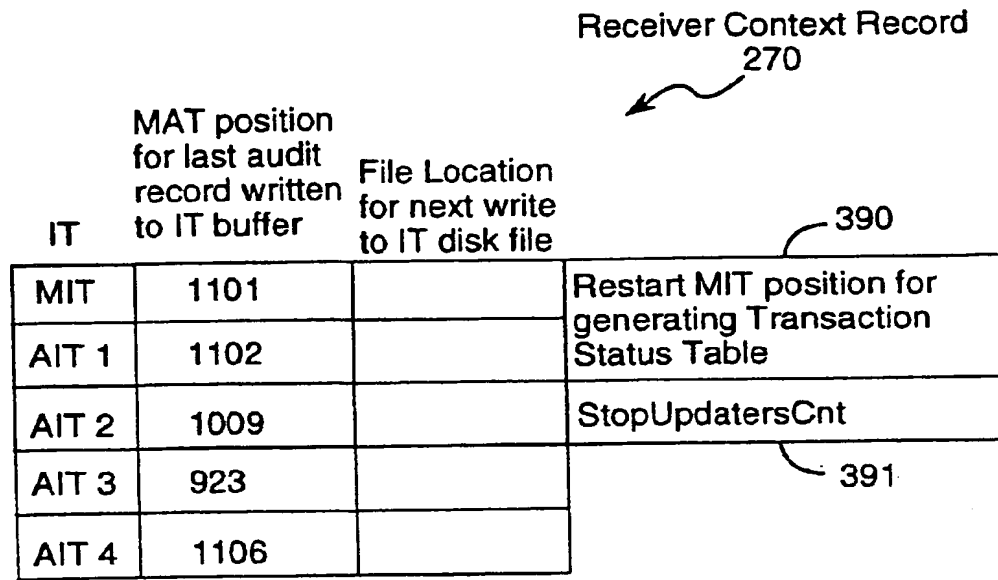


FIGURE 7A

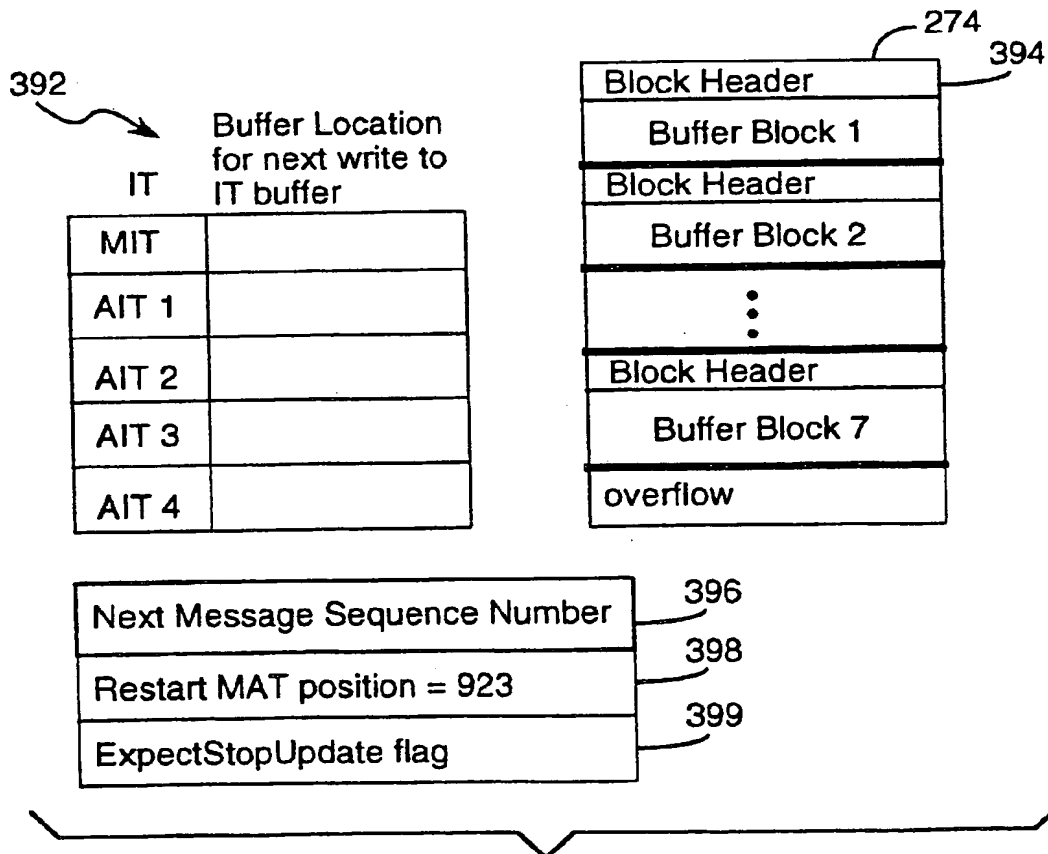


FIGURE 7B

Updater Status Table
400

Updater	IT File	Current StopPoint Location in IT File	Last Reported MIT position	Last Reported IT file position	Status
1	MIT				
2	MIT				
3	MIT				
4	AIT1				
5	AIT1				
⋮	⋮	⋮	⋮	⋮	

402
404
406
408
409

FIGURE 7C

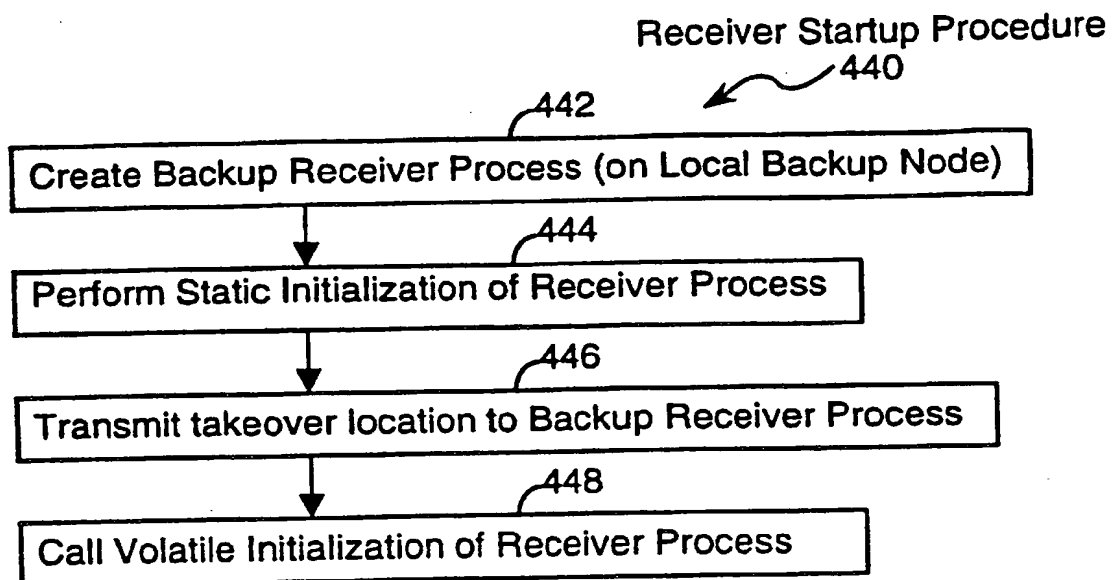
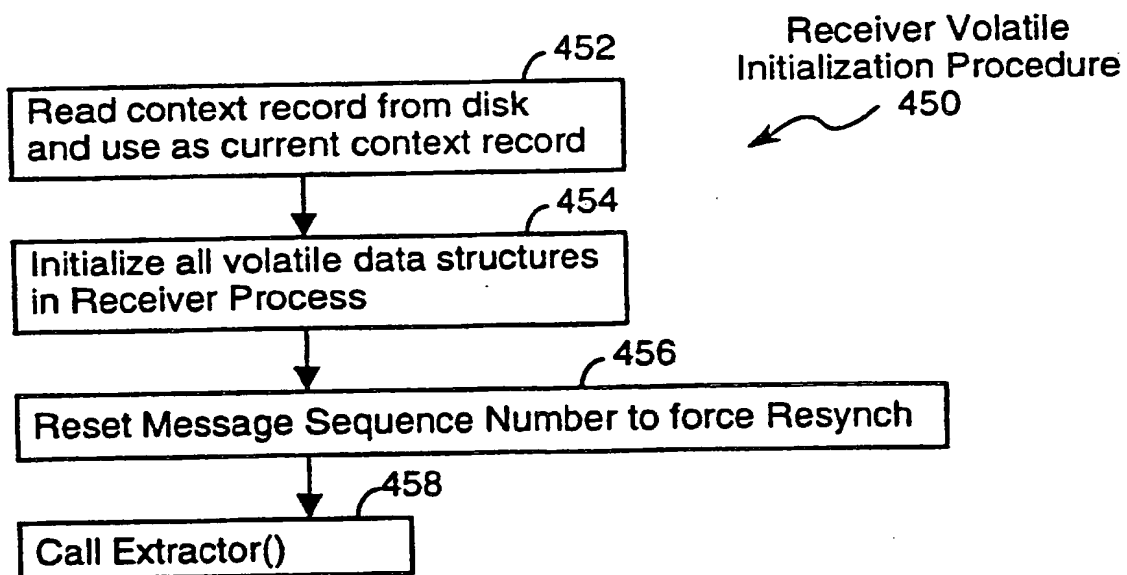
Transaction Status Table
(TST) 144

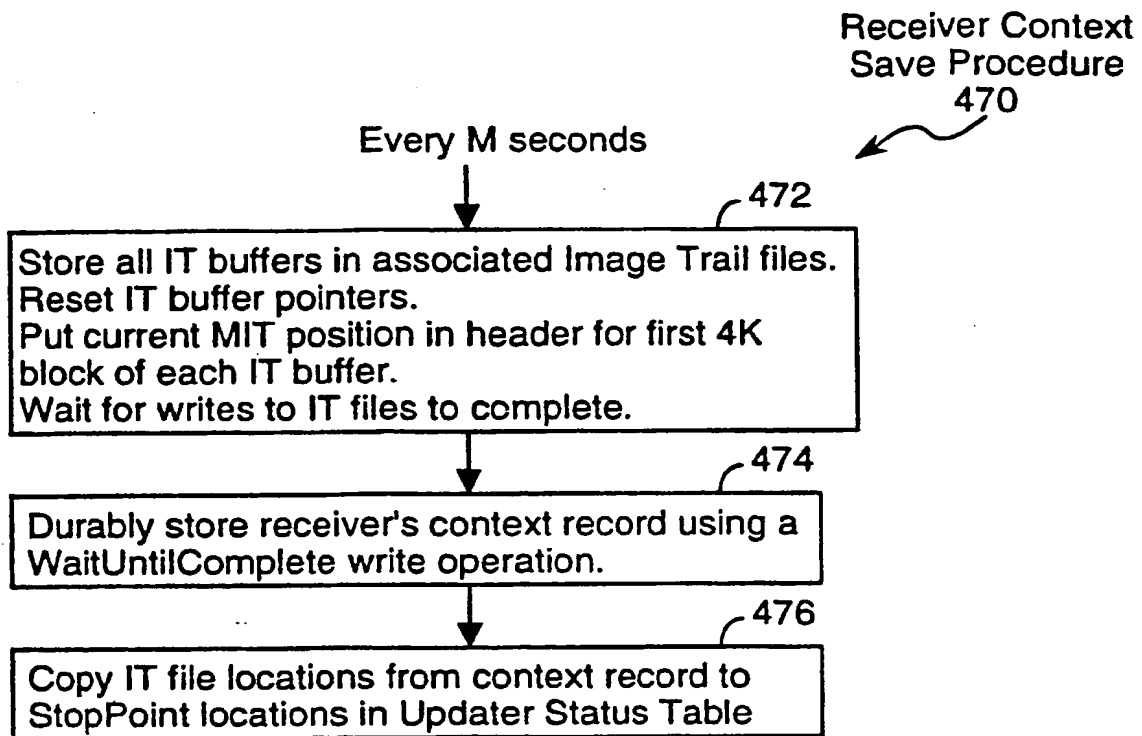
MIT position for
updating TST
414

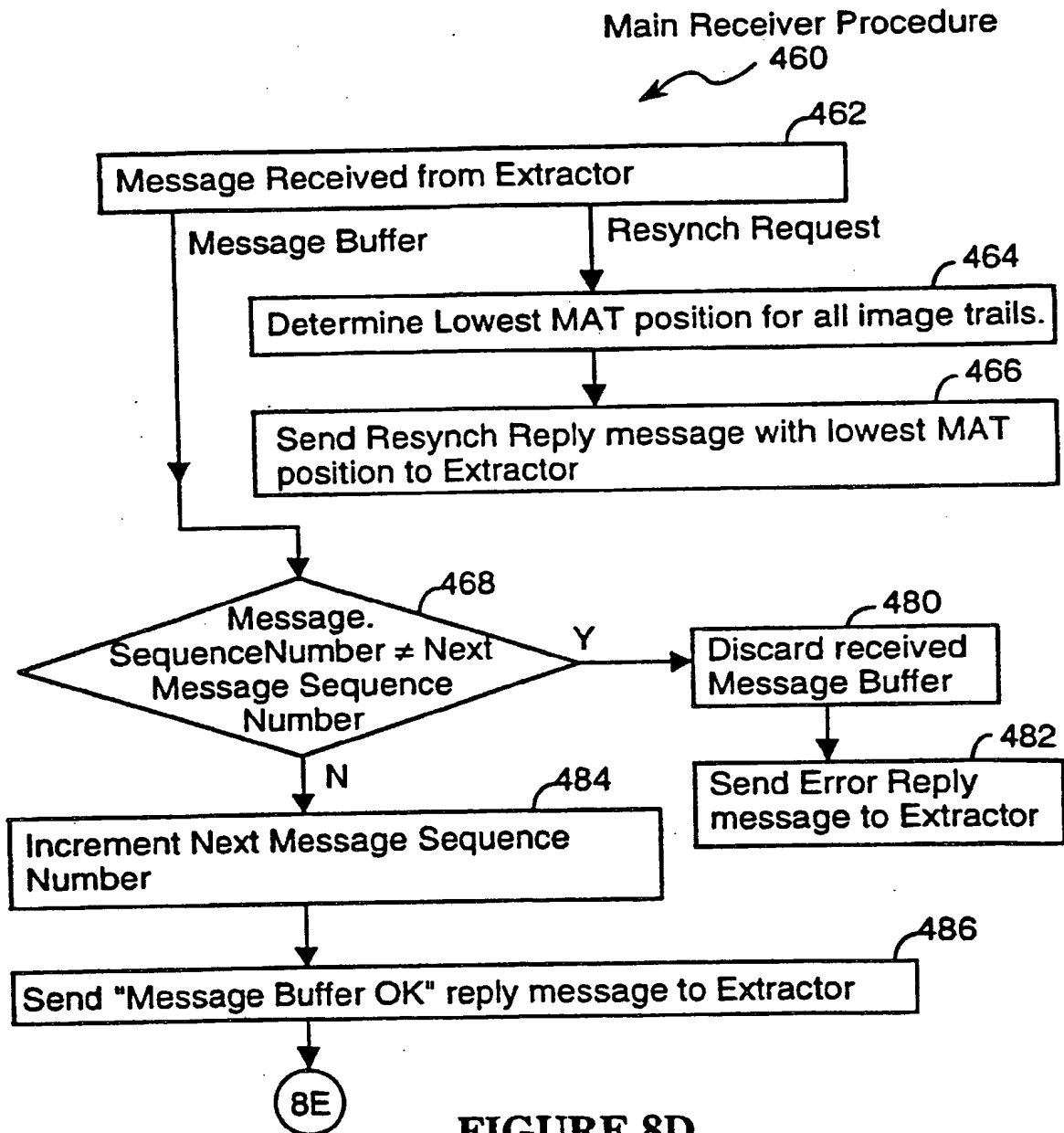
Trans. ID	Status (Commit or Abort)	MIT Position
S100003		
S100004		
S100006		
⋮	⋮	
S100256		

410
412
413

FIGURE 7D

**FIGURE 8A****FIGURE 8B**

**FIGURE 8C**



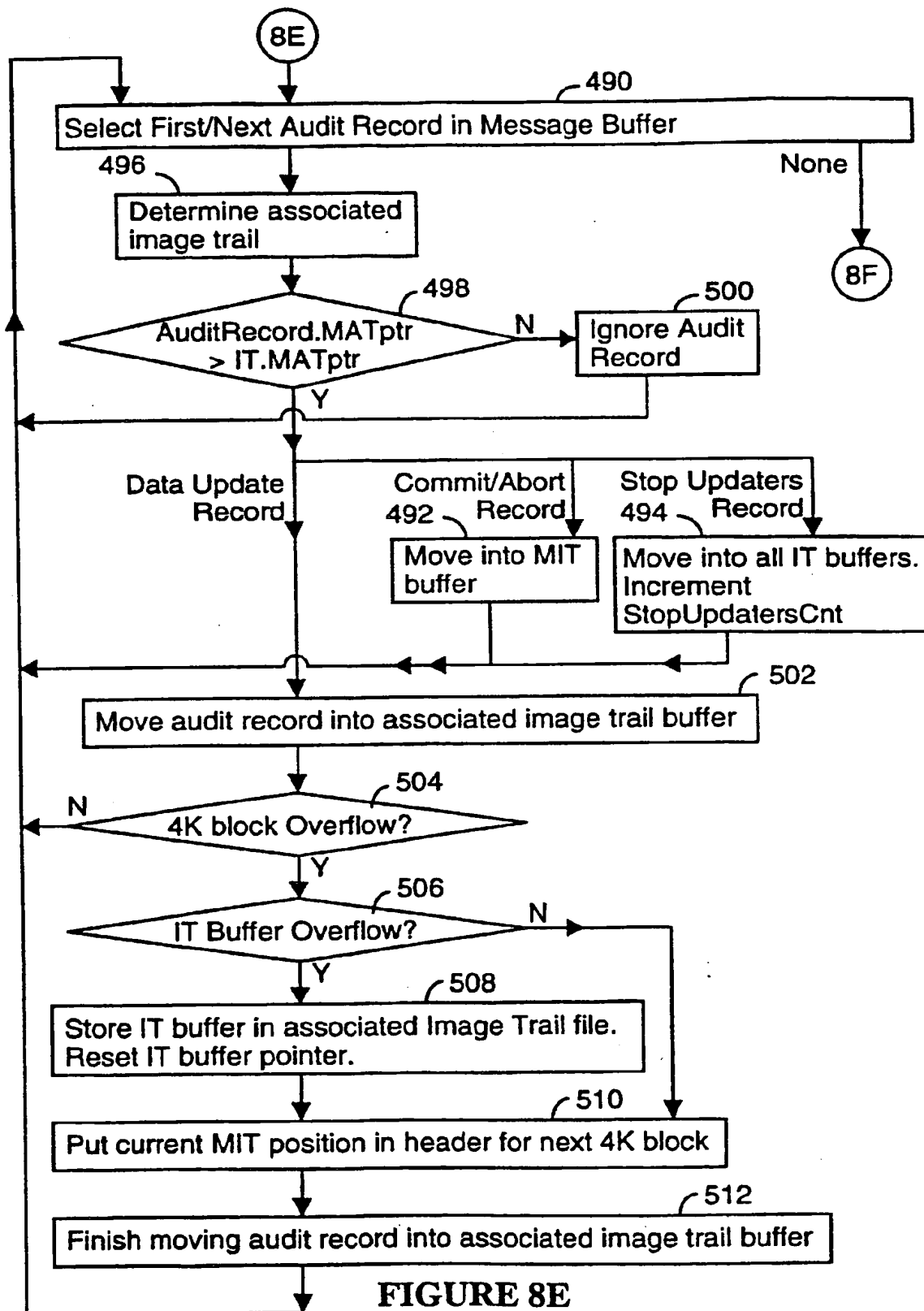


FIGURE 8E

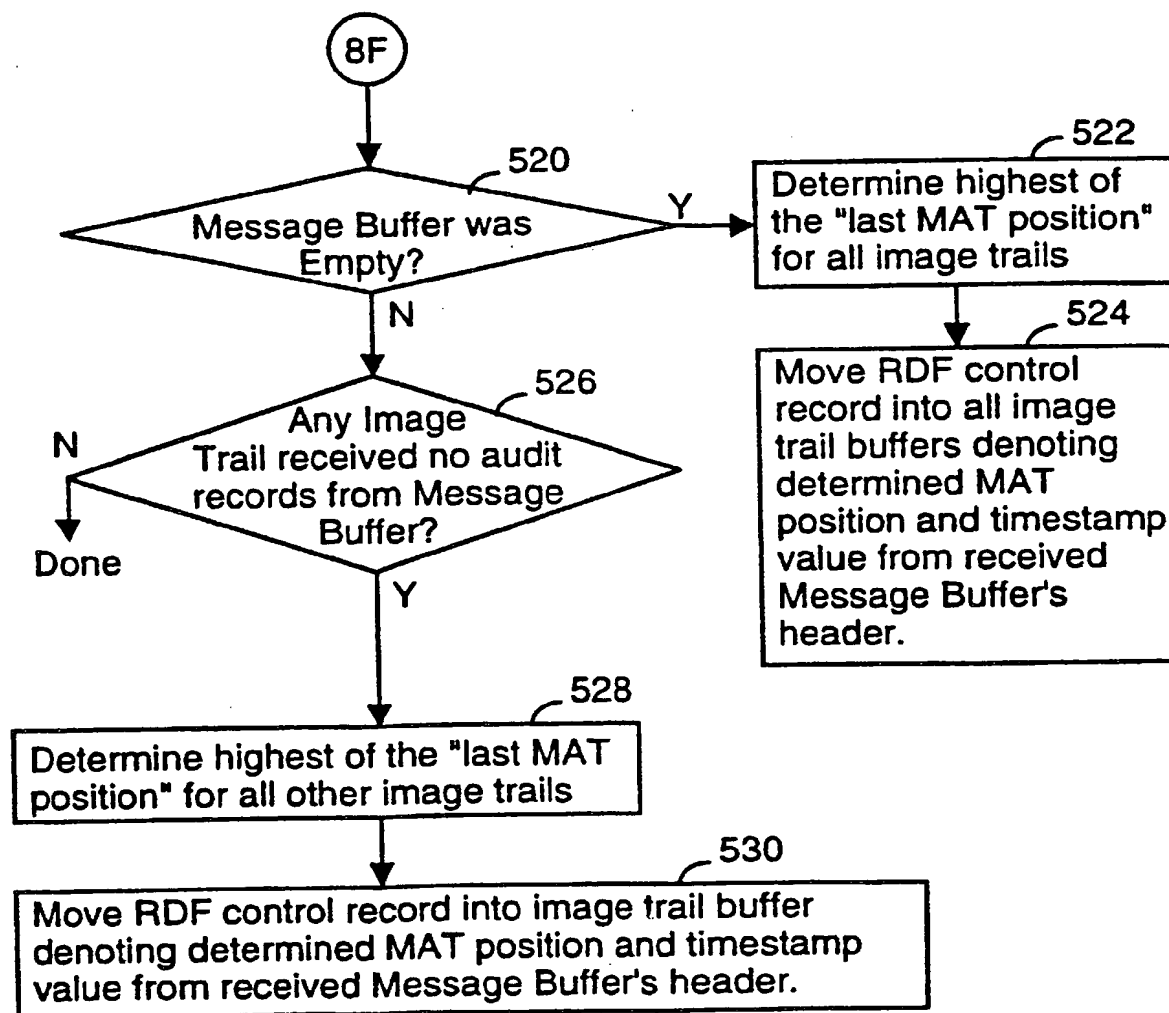


FIGURE 8F

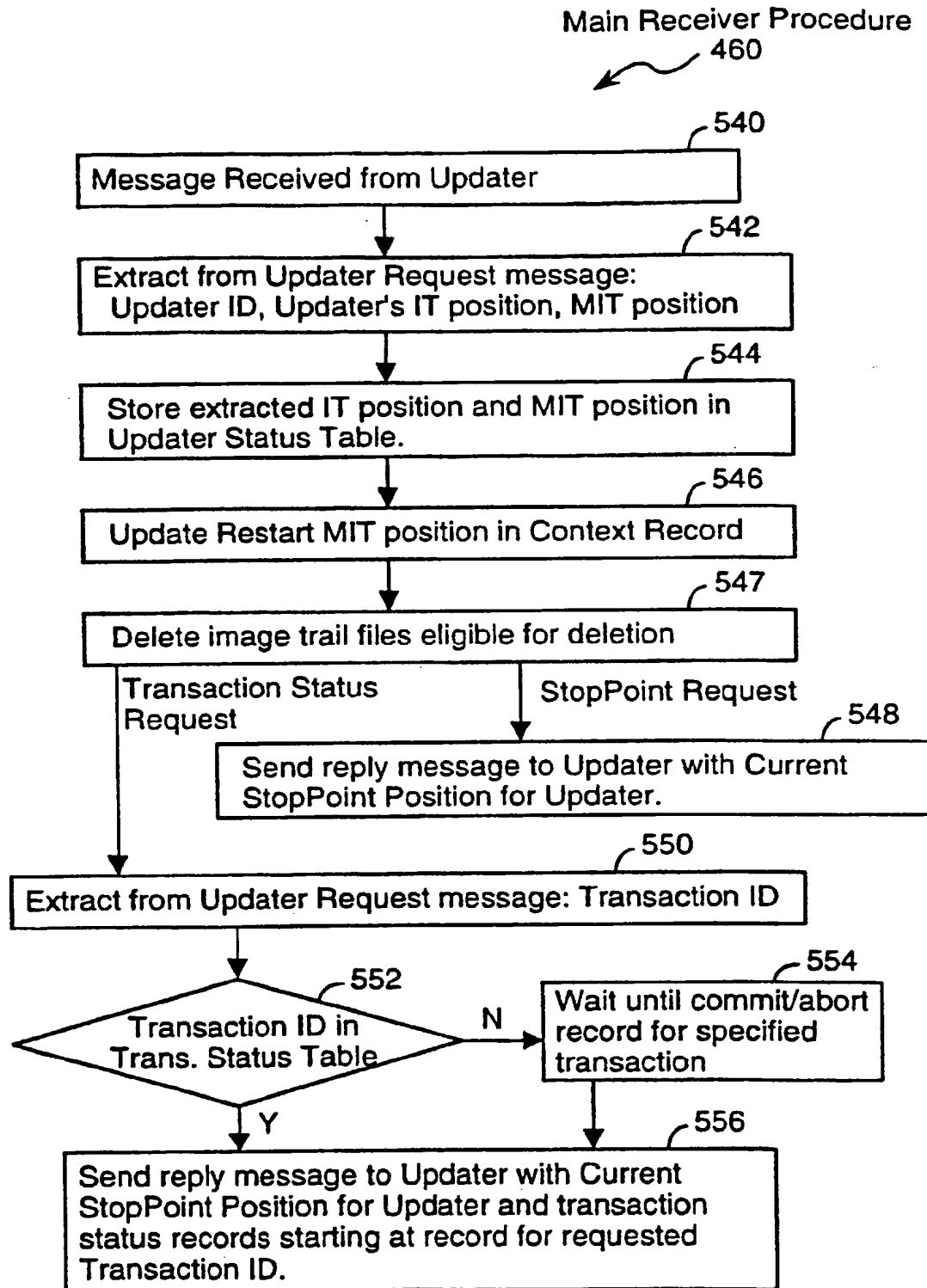


FIGURE 8G

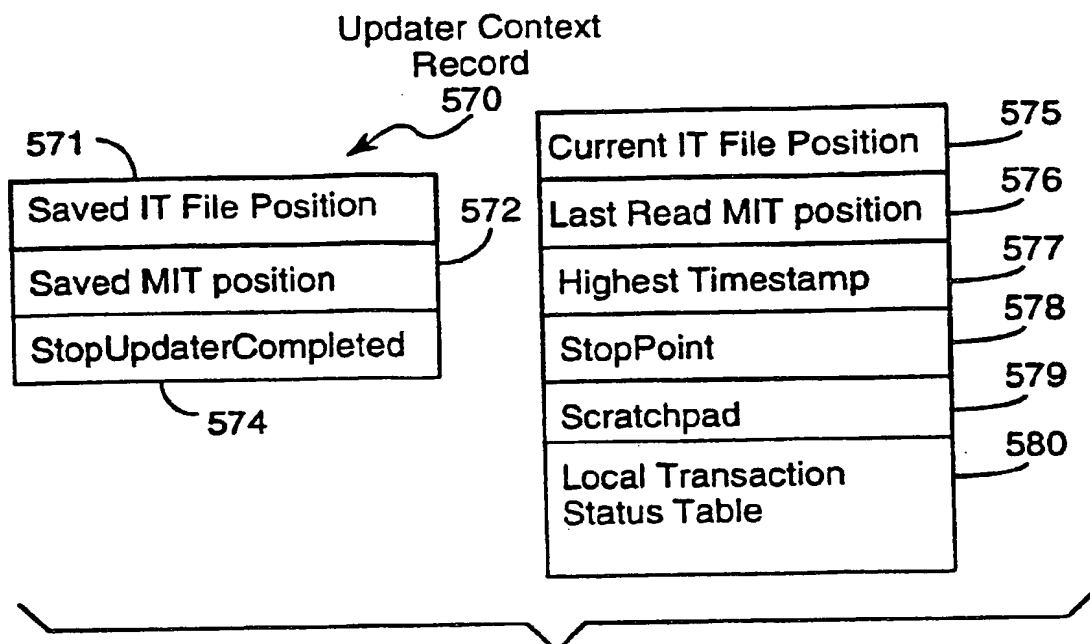


FIGURE 9

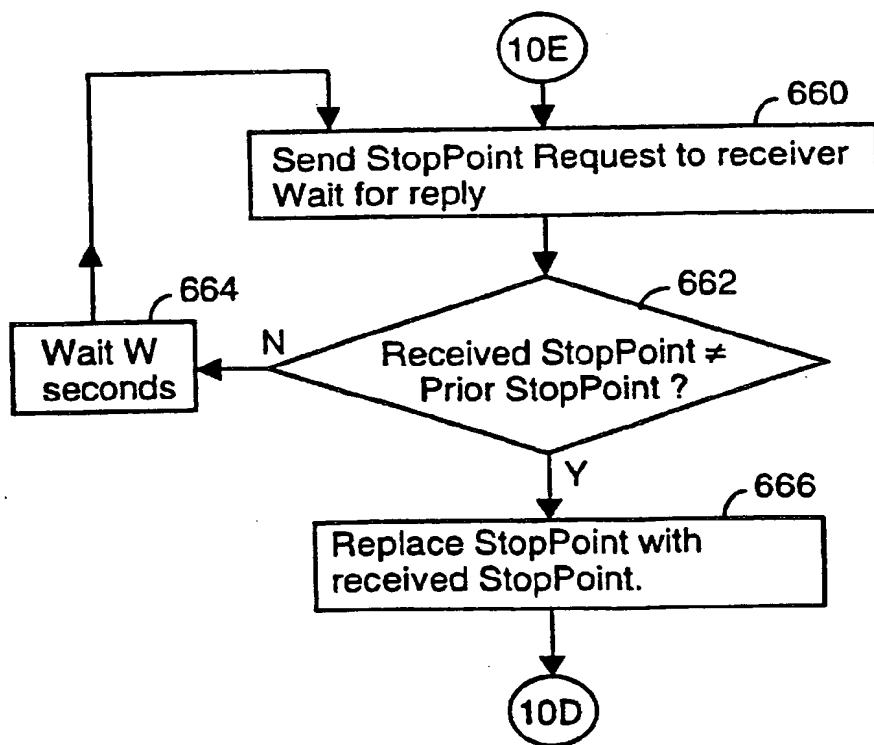


FIGURE 10E

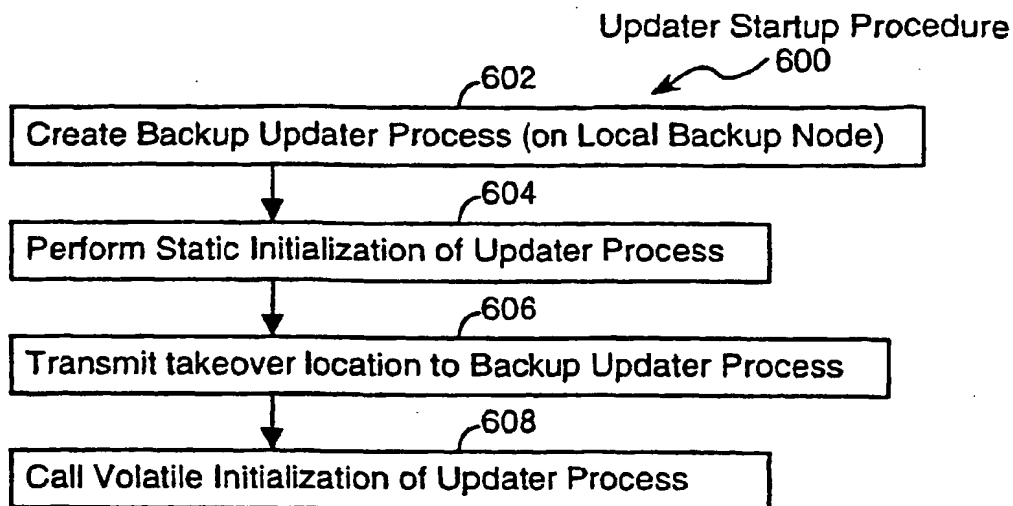


FIGURE 10A

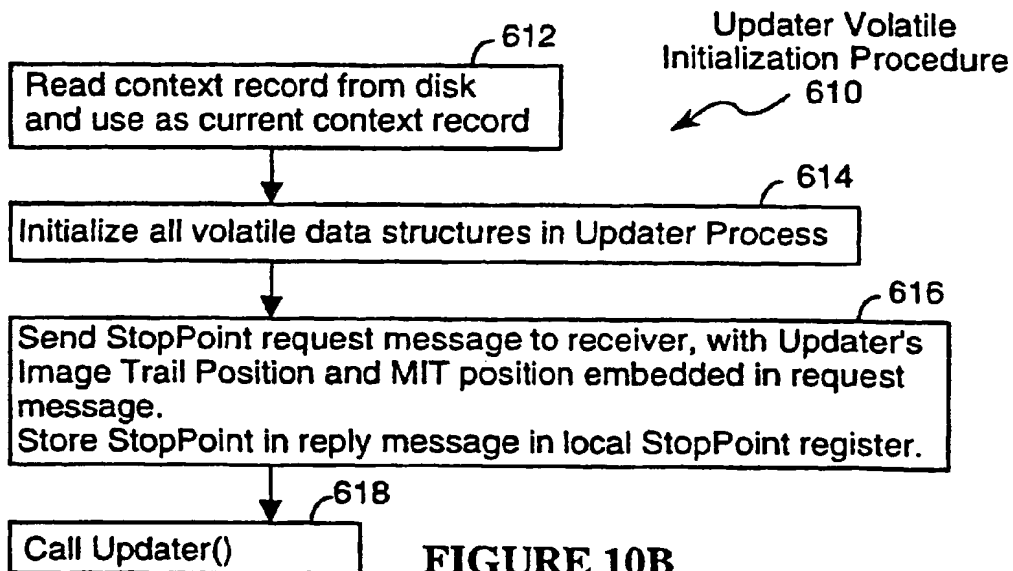


FIGURE 10B

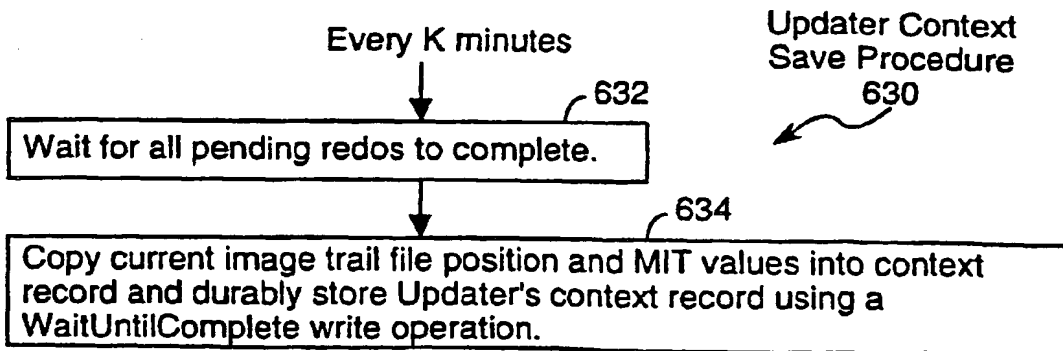


FIGURE 10C

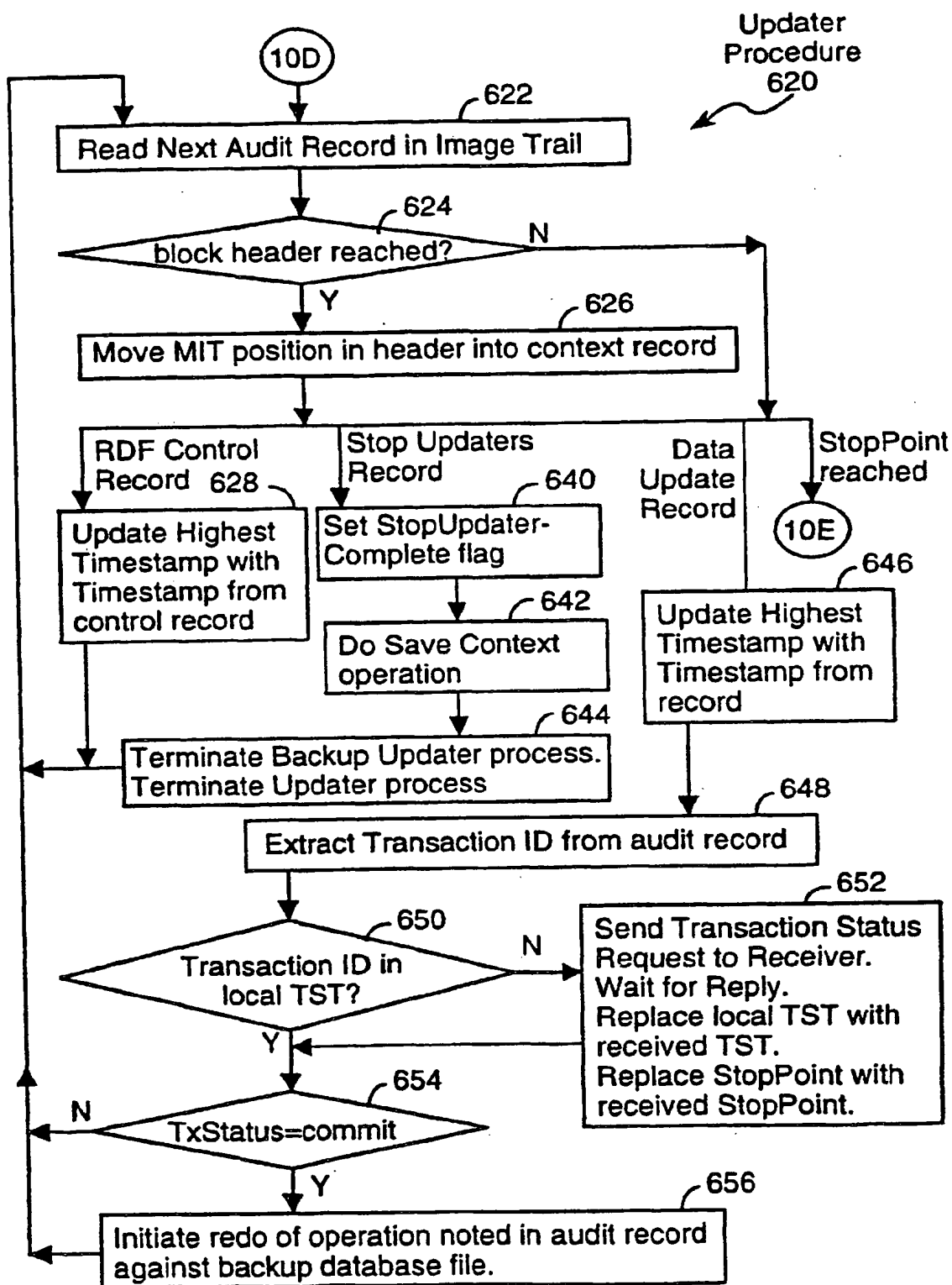


FIGURE 10D

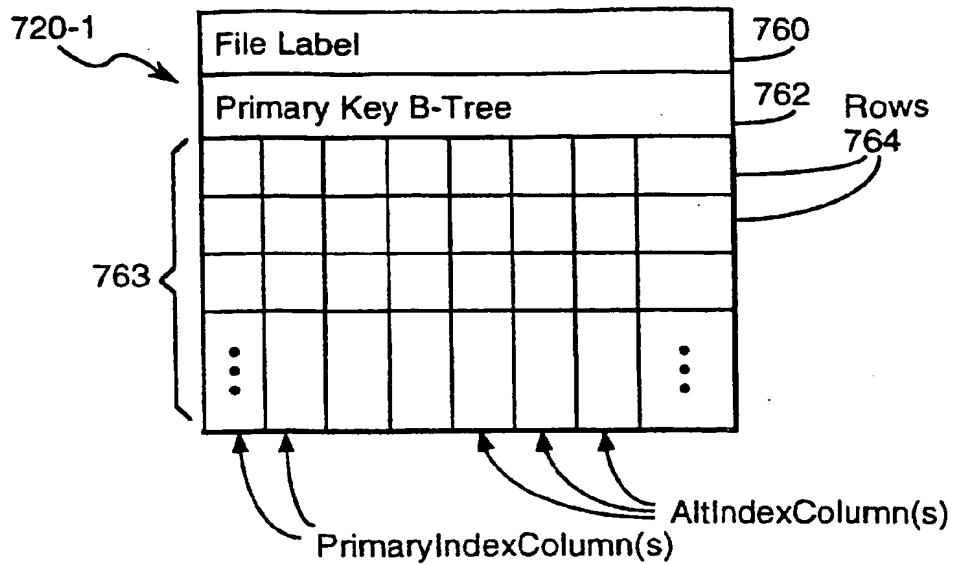


FIGURE 11A

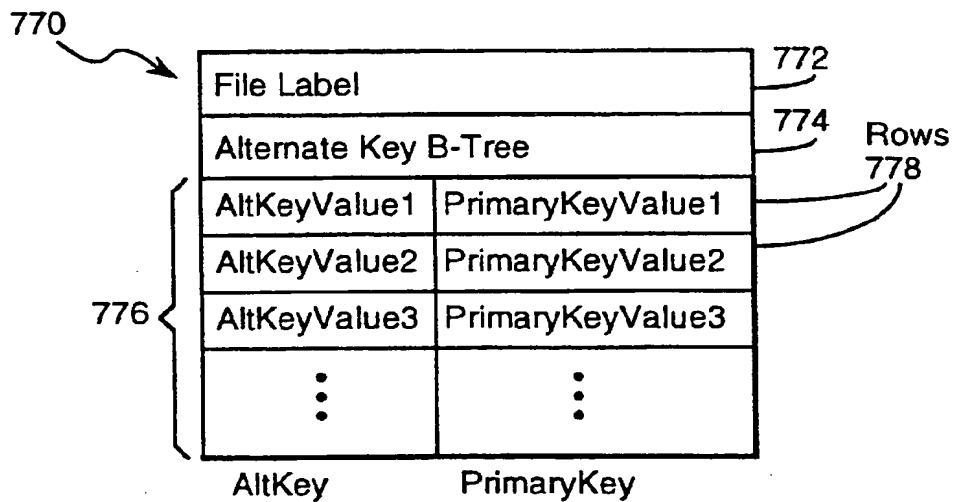


FIGURE 11B

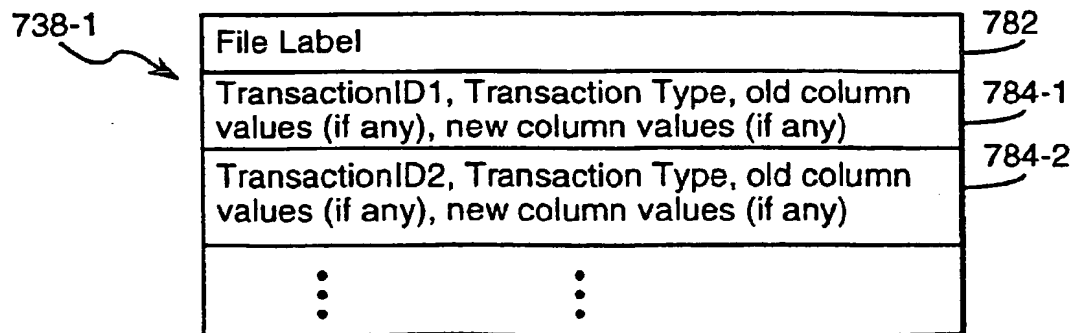


FIGURE 11C

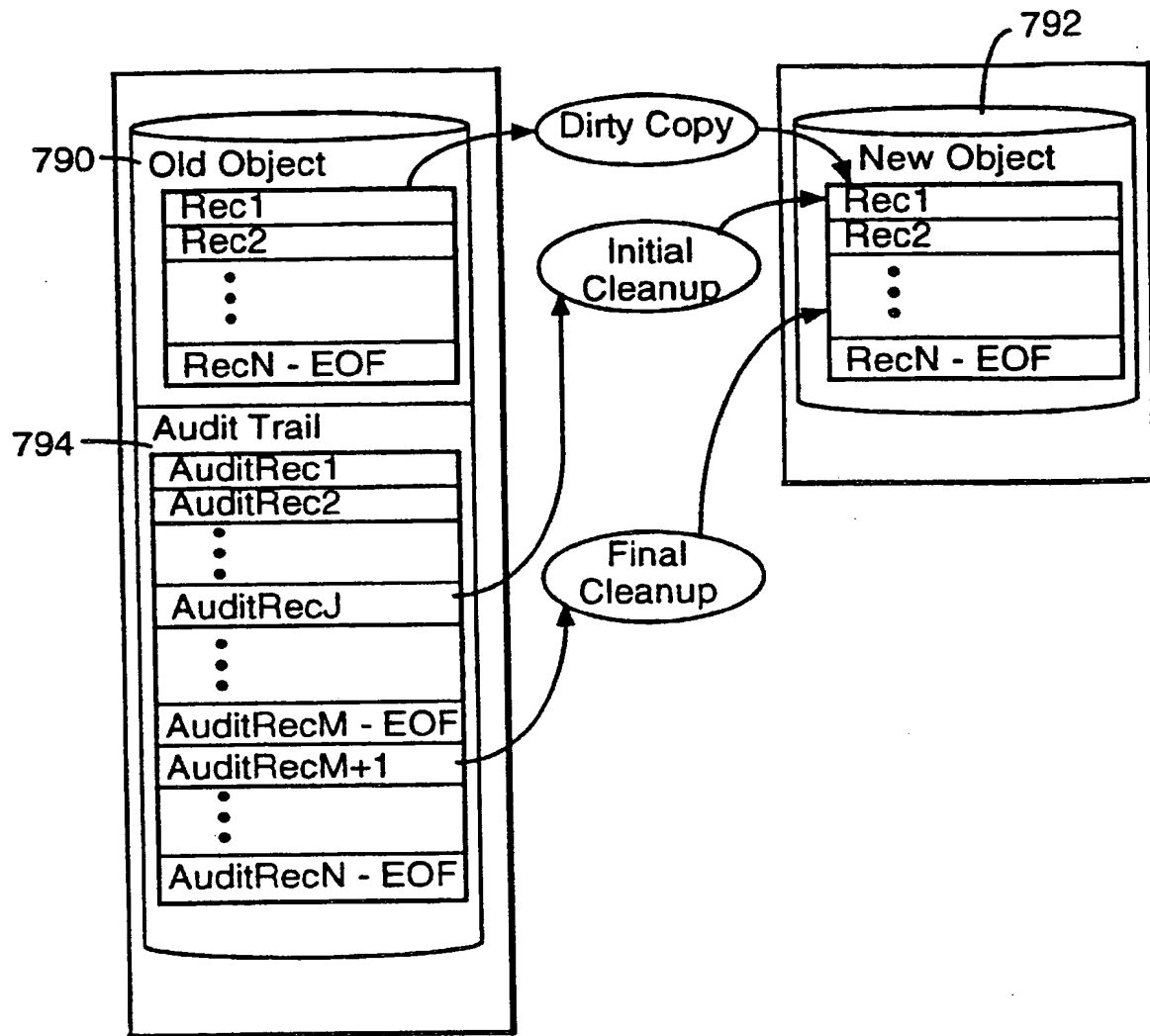


FIGURE 12

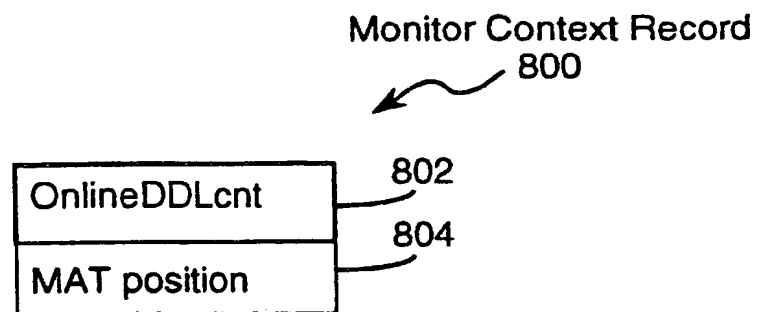


FIGURE 13

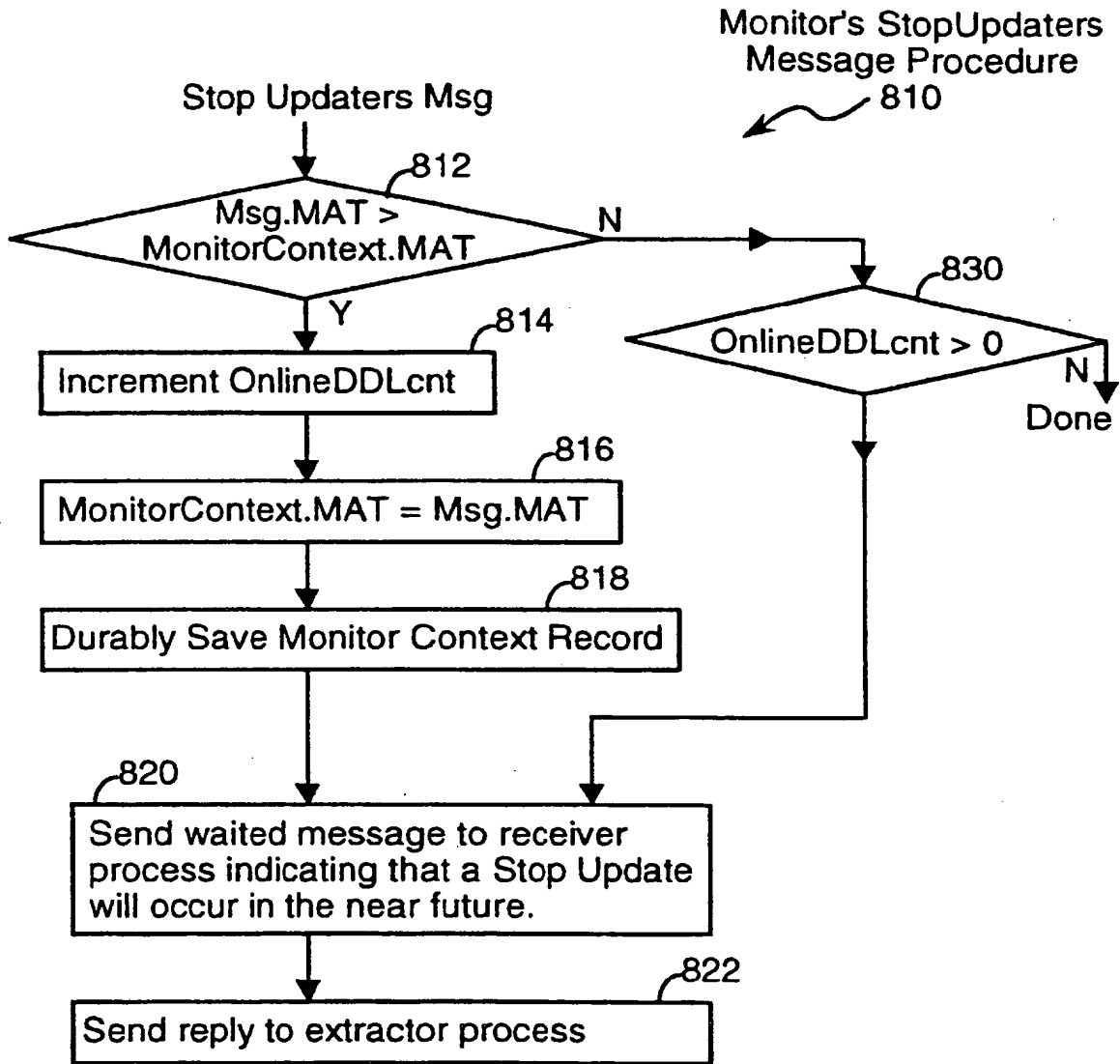


FIGURE 14

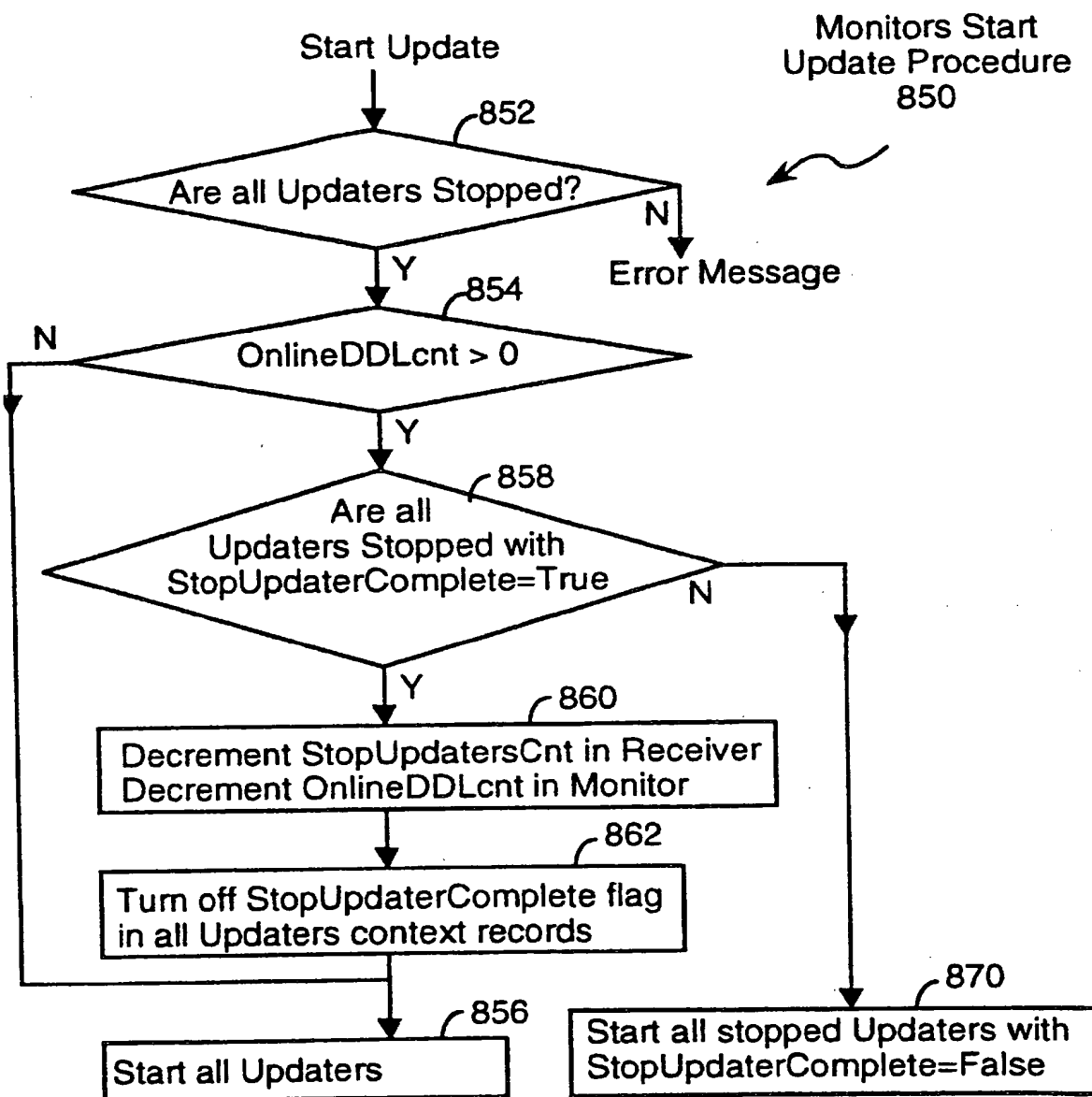


FIGURE 15

Receiver's Procedure for
Processing a Received
Stop Update Message

880

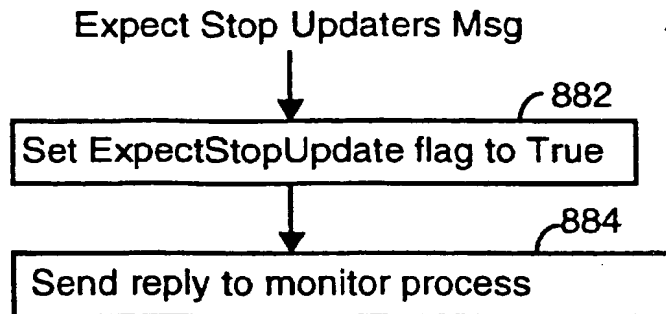


FIGURE 16A

Receiver's Updater
Close Message
Procedure

890

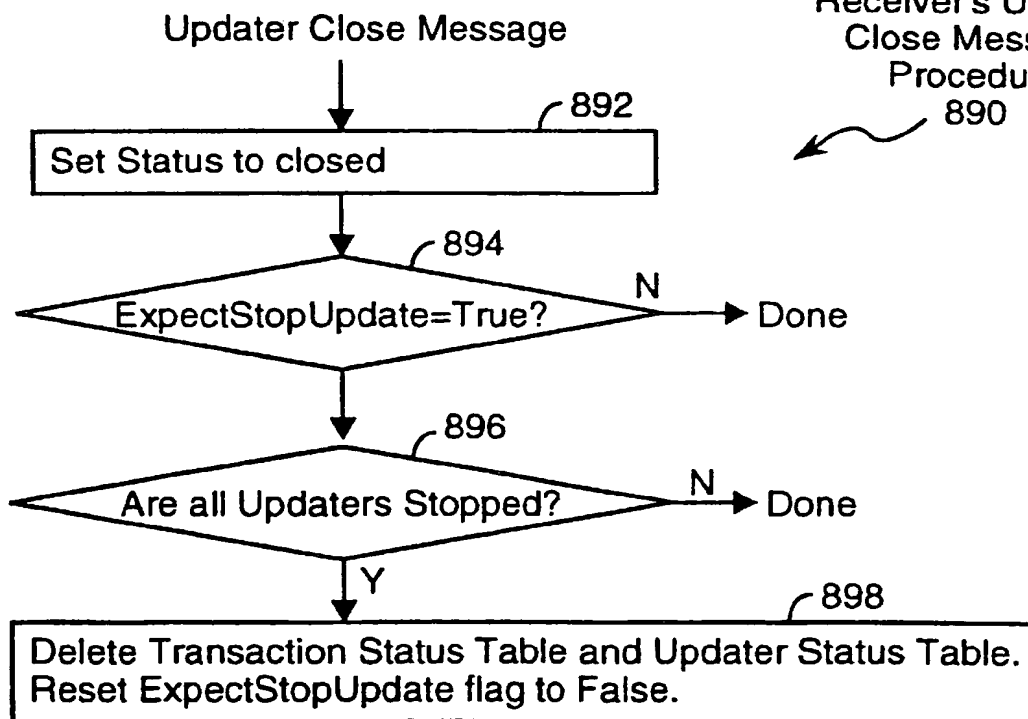


FIGURE 16B



European Patent
Office

EUROPEAN SEARCH REPORT

Application Number
EP 96 30 0437

DOCUMENTS CONSIDERED TO BE RELEVANT			
Category	Citation of document with indication, where appropriate, of relevant passages	Relevant to claim	CLASSIFICATION OF THE APPLICATION (Int. CL.6)
A	EP-A-0 600 457 (IBM) 8 June 1994 * abstract * * page 2, column 1, line 1 - page 3, column 3, line 32; figures 1-8 *	1,2	G06F17/30
A	US-A-4 875 159 (CARY RICHARD W ET AL) 17 October 1989 * abstract; figure 1 *	1,2	
			TECHNICAL FIELDS SEARCHED (Int. CL.6)
			G06F
The present search report has been drawn up for all claims			
Place of search THE HAGUE		Date of completion of the search 7 May 1996	Examiner Katerbau, R
CATEGORY OF CITED DOCUMENTS		T : theory or principle underlying the invention E : earlier patent document, but published on, or after the filing date D : document cited in the application I : document cited for other reasons & : member of the same patent family, corresponding document	
X : particularly relevant if taken alone V : particularly relevant if combined with another document of the same category A : technological background O : non-written disclosure P : intermediate document			

EPO FORM 1501 01.92 (P04C01)